

Are we biased on bias? Characterizing social bias research in the ACL community

Anonymous EMNLP submission

Abstract

Recent events in business, politics and society have shed light on the importance and potential dangers of Natural Language processing (NLP) in the real world. NLP applications have gained unprecedented popularity not just among scientist and practitioners, but also the general public. As we develop new methodologies and curate new benchmarks and datasets it is more important than ever to consider the implications and societal impact of our work. In this paper, we characterize the landscape of societal bias research within the ACL community and provide a quantitative and qualitative survey by analyzing an categorized corpus of 348 papers. More specifically, we present a definition of social bias based on ethical principals and investigate (i) types of bias, (ii) languages, and (iii) type of paper. We find that there is significantly more work on gender biases and English than other languages. Finally, we discuss the possible causes behind our findings and provide pointers to future opportunities.

1 Introduction

Traditionally, the NLP community has focused on ethical debates around privacy (Hovy and Spruit, 2016) ensuring that data is anonymized appropriately. More recently, there has been increased awareness that NLP research has a direct impact on peoples' lives (Mayfield et al., 2019; Bender and Friedman, 2018). For example, summarization systems can amplify misinformation (Smiley et al., 2017), and sentiment analysis (SA) systems can assign more negative sentiments/emotions based on race and/or gender (Kiritchenko and Mohamad, 2018). While such research used to be more academic (Leidner and Plachouras, 2017), these concerns are having an increasing impact in industry (Schnoebelen, 2017; Jin et al., 2021b) with consequences for users (Prabhumoye et al., 2021). It is well known that language data encodes demographics and biases (Bender and Friedman, 2018). There

is a risk that using such data can disclose inappropriate information about particular individuals, as well as undesirable attitudes towards individuals and groups (Hovy and Spruit, 2016; Eckert and Rickford, 2001) and social hierarchies (Blodgett et al., 2020). There are also concerns that systems based on inappropriate data are likely to repeat such biases, and may even amplify them (Bender and Friedman, 2018). In this paper, we survey 348 papers collected from the ACL anthology that focus on social bias and ethics in NLP research. We make three kinds of contributions, where (i) we present a working definition of social bias from a philosophy perspective, (ii) quantify our findings by annotating our corpus of papers and (iii) provide a discussion and pointers for possible future research directions. Through a quantitative analysis of current trends we attempt to answer the following questions:

- What kind of social biases is the ACL community concerned with?
- How many languages is bias studied in?
- What types of papers are present?

2 Related previous surveys

There is a considerable literature on social biases in NLP. Much of this work provides guidelines and/or recommendations. One of the first position papers on the topic outlined the need for ethical considerations that go beyond privacy concerns for users, and focus on the social impact of experiments and applications on individuals and (minority) groups Hovy and Spruit (2016).

Surveys on social bias emphasize a variety of different aspects, such as embedding representations, data collection and annotation, downstream task performance, metrics and limiting negative impact (Garrido-Muñoz et al., 2021; Mohammad, 2022b; Schnoebelen, 2017). Work by Bender and Friedman (2018); Hovy and Prabhumoye (2021) outline the concept of data statements and sources

of bias respectively to aid the research design process. Other research has reviewed how to mitigate bias (Chandrabose et al., 2021; Meade et al., 2022; Balkir et al., 2022), how to teach bias, ethics and privacy to students (Bender et al., 2020; Friedrich and Zesch, 2021), evaluate existing metrics (Czarnowska et al., 2021; Delobelle et al., 2022), handle challenges presented by new laws (e.g., GDPR) (Lewis et al., 2017) and apply existing principles from ethics and privacy to NLP (De Jong et al., 2018; Leidner and Plachouras, 2017; Prabhu-moye et al., 2021). Our work follows the precedent established by Blodgett et al. (2020), who used keywords to select papers from the ACL anthology, and then enlarged the sample by following citations to other popular venues (eg AAAI, ICML etc.). Similarly, we also align with Field et al. (2021) who solely focus on papers published at ACL venues but draw on conclusions from NLP papers published elsewhere. Relying on keywords, of course, introduces possibilities for false positives and false negatives.

Scope of this survey In this work, we solely focus on papers published in the ACL anthology to gain a better insight into current trends and popular research approaches for social bias. One limitation of such an approach is that seminal work published at other venues is not reviewed here. In Appendix 7, we provide a table that shows each paper reviewed in this survey.

3 Defining Bias

The dictionary definition of bias is “an inclination or prejudice for or against one person or group, especially in a way considered to be unfair” (Stevenson, 2010). Based on this, defining bias involves an interaction among different components: (i) individuals or categories that determine a group, (ii) attitudes towards this group, and (iii) assessment of this attitude in relation to fairness. Bias comes into existence when a specific attitude is formed, which may or may not be fair. For example, in investigating a disease that is more prevalent among women, use of *gender* is a relevant variable and, in itself, does not entail a differential attitude. Once a differential attitude is formed however, bias comes into existence and the attitude may be negative or positive. By definition, if a distinctly positive attitude is formed about one variable (e.g., *gender*), it entails a less favorable attitude about the other *gender* categories. This does not mean that all bias

is necessarily unfair, there are multiple theories and definitions of fairness that are formulated and analyzed in-depth in political philosophy (Lamont and Favor, 2004). The formal principle of equality formulated by Aristotle (Ameriks and Clarke, 2000) states that equals must be treated equally, which is often referred to as ‘the fairness ideal’, but it is neither a prevailing definition nor a useful one in practice. Without identifying relevant features, such a definition would not prevent categories such as *race* or *wealth* to be used as variables for differential treatment. However, a fairness approach (Lamont and Favor, 2004) based on equal opportunity might require a differential attitude (i.e., bias) towards a certain category in order to ‘level the playing field’. For example, if women are routinely given worse performance reviews and lower pay for successfully completing the same tasks as men, then there is differential treatment. Meaning women who are as successful as men cannot have the same opportunities. According to this understanding of fairness, a bias towards women would be fair. It is also worth noting that positive or negative discrimination is distinct from bias. While discrimination is about treatment, bias is about attitude. In other words, bias may lead to discrimination (Mateo and Williams, 2020). In this context, the ethical issue about bias is tied to differential treatment of descriptive categories resulting in unfair outcomes. Identifying and dealing with ethical concerns related to bias, must necessarily involve identifying the descriptive categories and the biases against those categories as well as examining whether the said bias is unfair, according to the relevant definition of fairness.

4 Methodology

Following the standard practice mentioned above, we searched the ACL anthology¹ in September 2022 for relevant titles and abstracts using the keywords: *ethical, ethics, fairness, fair, bias, social, society, societal, social good*. For papers, published after September 2022 we manually screened all conference proceedings for the same keywords. One limitation of relying on a keyword search alone is that we might miss any work that refers to a bias directly in the title, for example ‘fatphobia detection in online forums’.

¹<https://github.com/shauryr/ACL-anthology-corpus>

Filtering Strategy A total of 1,437 papers were returned by the search; 523 papers were retained after a manual screening of titles and abstracts. We removed duplicates, as well as work not related to bias and/or ethics in NLP. Then we downloaded full papers, and filtered out papers if: (i) the contribution was a talk, demonstration, abstract or keynote, (ii) “bias” was used in the machine learning sense, or (iii) the paper did not focus on social bias. This process produced a collection of 348 papers.

Categorization process We identify trends in our corpus by empirically determining a set of five categories, where we fully review each paper manually. We make our corpus alongside with its corresponding categories and labels publicly available². We focus on four elements for each paper to identify trends, where we identify language investigated, type of bias analyzed and what kind of paper is introduced and which NLP area it belong to. For the type of paper analysis we utilize the authors description of contributions to split the papers in the following categories:

- **Method:** In this category of papers, the main focus of the work is to contribute a new method, which includes but is not limited to novel ways to measure or mitigate social bias.
- **Analysis** Papers in this category, examine existing datasets, benchmarks, language models, NLP systems or embeddings for bias using social science methods, statistics and mixed methods approaches. For example, authors who have conducted research in this category have explicitly stated that they conduct an analysis or outline a mixed method approach.
- **Surveys and Position papers:** This category of papers includes surveys, guides, tutorials, reviews and position papers.
- **Dataset, Benchmarks or Resources:** Paper in this category propose new datasets, benchmarks, lexicons, challenge sets and often include some preliminary analysis of the new data either collected through crowd-sourcing.
- **Datasets, Benchmarks and Methods:** This is a combination of papers that focus on both introducing a new resource (e.g.: dataset or benchmark) in addition to a new methodology.

²Link-added-upon-publication

Year	Papers	Year	Papers
2010-2016	11	2020	68
2017	16	2021	96
2018	10	2022	79
2019	46	2023	22
Totals	83	Totals	265

Table 1: 76,15% of the 348 papers are from 2020-2023.

5 Empirical Findings

The 348 papers were published between 2010 and 2023. Table 1 shows that there has been considerable growth in interest in bias research.

Type of Paper Based on the criteria outlined above, we have found that the majority of papers introduced are *Method* papers (57.75%), followed by *Datasets and Benchmarks* (20.40%), *Surveys* (15.22%), *Analysis* (4.31%) and combined *Datasets, Benchmarks and Methods* (2.29%). We take a closer look the majority category *Methods*, where we assigned *Method* papers into different categories based on the contribution of the paper (see Table 2).

Methods We distinguish between bias detection and measurement, as detection does not necessarily measure or remove any kind of bias. We found similarly to (Blodgett et al., 2020) that many papers propose a combination of techniques, hence why we have decided to merge such approaches into the category ‘*Measurement and Mitigation*’. Many works in Debiasing and Mitigation apply or extend methods such as WEAT (Caliskan et al., 2017) or HARD DB (Bolukbasi et al., 2016) to a specific benchmark or dataset. One negative side effect of this could be that for example in *gender bias* Measurement/Mitigation, there is no evidence that HardDB can or should be applied to languages with grammatical gender (Sun et al., 2019). There has also been criticism of removing bias (Caliskan et al., 2017), where removing bias (i) only changes how an application or algorithm understands the world but not how it applies the knowledge gained from understanding (‘fairness through blindness’), (ii) could harm meaning and accuracy and (iii) bias can (unintentionally) be introduced through other avenues during the design process. Therefore, simply removing bias is not enough (Chandrabose et al., 2021) and developing new methods requires consistent reflection as bias in NLP systems is never fully inescapable (Waseem et al., 2021). Work by (Hovy, 2015; Kiritchenko and Mohammad, 2018) has found that the some

information bias mitigation techniques can be beneficial in improving performance in downstream tasks. The methods in the *Miscellaneous* category take a different approach to working on bias. For example (Fisher et al., 2020) in including bias sensitive attributes by defining a whitelist of triples for uncontroversial cases. Arguably, one drawback here is that it is up to the person whitelisting to decide what is not controversial and what is. Similarly, (Touileb et al., 2021) and (Wang et al., 2017) use *gender* and *linguistic bias* respectively to improve classification results. Touileb et al. (2021) first show how female critics disproportionately give lower ratings to female authors, where removing metadata may have the opposite effect in that it does not help traditionally underrepresented groups in a specific domain. At the same time, many of the methods looking at measuring or removing bias complicate the data and tasks at hand and can lead to the development of systems that are not reliable when used in a more complex context (Talat et al., 2022b). This also applies to the predominant use of intrinsic metrics in bias measurement. These metrics may shed more light on how much bias exists in a dataset/LM, but does not necessarily correlate with performance on downstream tasks and therefore does not show the true harm of bias (Orgad and Belinkov, 2022). Thus, we may run the risk of developing methods for each new dataset or benchmark and missing out on crucial information that shows how bias affects different downstream tasks in different ways. However, documenting and measuring bias in a systematic way is crucial to understanding what harms can be caused in real life situations, so that preventive methods can be developed (Dev et al., 2021b). Current approaches in mitigation and/or measurement methods are evaluated on a variety of NLP areas, including Language Models (35.74%), Classification (22.85%), NLG (14.76%) and NLI (3.33%). It is unsurprising that much attention has been paid to embedding representations that are trained on large amounts of text (Kiritchenko and Mohammad, 2018; Mayfield et al., 2019; Talat et al., 2022b). This has the benefit of bias methods being more widely applicable, but it also means that there are distinct limitations when a method is tied to a specific architecture rather than the task/benchmark itself. It means bias measures are no longer comparable in relation to other benchmarks and bias can be introduced at any stage of an NLP system design as it de-

Type	Papers	%
Measurement	98	46.66
Debiasing / Mitigation	52	24.76
Combinations of above	32	15.23
Detection	19	9.04
Generation	5	2.38
Miscellaneous	4	1.9
Total	210	100.0

Table 2: Empirical taxonomy of methods. depends on where and how the final LLM is applied and to which community (Talat et al., 2022b). An important trade-off to consider is the balance between generalizable and context-sensitive methods to measure bias in downstream tasks. There are also other areas that have done work on bias but are not represented as well in this survey, which include but are not limited to Speech Recognition (Kwako et al., 2022; Savoldi et al., 2022), Multi-Modal NLP (Chen et al., 2020a; Srinivasan and Bisk, 2022a) and Information Extraction (Li et al., 2022b; Sun and Peng, 2021).

Social biases In Figure 1 we show the types of biases investigated, where we include all social biases that occur more than once (see diagonal of the matrix). Furthermore there are a small number of biases that only occur once and are not shown in the table, such as bias transfer hypothesis (Steed et al., 2022) or dialect bias (Tatman, 2017b). Furthermore, we included the category *multiple social biases*, where the paper does not explicitly list or describe the specific type of bias investigated (Ghosh et al., 2021; Ramponi and Tonelli, 2022; Mireshghallah and Berg-Kirkpatrick, 2021; Loukina et al., 2019). There also is the *intersectional bias* category, which shows how different elements of a person’s identity (e.g., gender, race and age) can either be a benefit or disadvantage and lead to compounded discrimination (Crenshaw, 2013). For example, recent work by Lalor et al. (2022) benchmarks a variety of NLP models on different downstream tasks for its performance on intersectional biases and (Cámara et al., 2022) introduce a framework for unisectional and intersectional analysis of sentiment analysis in a multilingual setting. Few works focus on biases other than *gender*, where (Davidson et al., 2019; Sap et al., 2019; Manzini et al., 2019) look at *racial bias* and (Hutchinson and Mitchell, 2019; Herold et al., 2022) investigate *disability bias*. Noticeable is that some biases are not investigated on their own, such as *age*, *religion*, *sexuality* or *profession*. We included *political/media bias* in this analysis, if the paper also looks at social bias. For example, debiasing claims that include

attitudes towards a group (e.g., *sexuality*). However, this type of work does not explicitly mention social biases when attitudes or characteristics of the targeted group are only implied. Dayanik and Padó (2020) looks media bias on a immigration dataset (MARDY) but does not mention implied social biases (e.g., *nationality* or *ethnicity*). The most frequently combined social biases are *gender* and *race*. The most frequently combined social biases are *gender* and *race*. In Table 3 we compute residuals between observed joints and predictions on margin, where highlighted in *green* are highly saturated areas and shades of *red* show less saturated areas.

Languages We show the languages used in each type of paper, excluding *Surveys and Position papers* in Figure 2. There are a total of 34 languages, however we leave out any languages that only occur once in the visualization. There are 11 languages not visualized, including *Farsi, Urdu, Wolof, Bengali, Armenian, Bengali, Inuktitut, Ukrainian, Hungarian, Indonesian* and *Lithuanian*. English, German, Spanish and Chinese are most commonly used either on their own or in combination with each other. The majority of all papers focus on a single language at a time. Furthermore, the vast majority of LLMs are monolingual and do not encode the cultural variety that naturally occurs within one language, for example non-standard English varieties (Talat et al., 2022b). Therefore it is important to not only document the type of bias investigated, but also contextualize bias within a language’s cultural context, understanding of said bias and document the language itself (*Bender Rule* (Bender, 2011)). Based on this collection, bias research is heavily biased towards western and Anglo-centric notions of bias and very few works focus on non-English benchmarks (Talat et al., 2022b; Hovy and Spruit, 2016). This proves extremely problematic when English benchmarks are automatically translated, but many of the biases do not hold true in non-Western cultural contexts. For example, gendered professions do not necessarily translate across every language or culture (Talat et al., 2022b) and many NLP systems trained on written English (e.g., Penn Tree Bank) do not perform well on non-standard English (Mayfield et al., 2019). From Table 4 we can see the residuals between observed joints showing a clear over-saturation (*green*) for specific combinations of languages (e.g: *English and German* or *English and Spanish*).

6 Discussion

Datasets and Benchmarks Previous work (Hovy and Spruit, 2016; Hovy, 2018; Talat et al., 2022b) outlined a number of reasons that may explain why there are so many papers focusing on the same datasets, benchmarks, biases, and languages. In the following section, we highlight some of the elements that may explain why there is an uneven distribution of work and resources.

- **Experimental setup** The majority of of work in bias research has focused on using intrinsic bias measurements (bias in internal model representations) and little attention has been paid to extrinsic metrics (Orgad and Belinkov, 2022; Delobelle et al., 2022). A very real consequence of this is that much work does not appropriately describe, contextualize and identify the potential harms that bias has in real-world scenarios (Blodgett et al., 2020). Also bias is inadvertently introduced in intrinsic metrics, where lexicons used to measure bias in one dataset produces very different results in another (Antoniak and Mimno, 2021). Similarly, (Goldfarb-Tarrant et al., 2021) have found that there is no correlation between intrinsic and extrinsic metrics. Another key problem is that often newly proposed datasets are linked to specific metrics, which makes it hard to draw conclusions from individual case studies as many results are not generalizable (Orgad and Belinkov, 2022). Another element impacting metrics is the composition of test data, where (Orgad and Belinkov, 2022) found that test sets often don’t contain balanced examples. However, most metrics are defined over a whole dataset and are therefore sensitive to its composition, which may lead to variability in results. Both metric choice and dataset composition can significantly change the results and conclusions drawn from a downstream task or dataset (Akyürek et al., 2022a). Therefore, it is important to (i) provide motivation for including / excluding a particular metric and describe how it impacts downstream performance, (ii) compare a metric across a variety of datasets and (iii) compare many metrics across individual datasets.
- **Funding** There are unintended consequences of research that can be traced to how research projects are funded (e.g., governments or mil-

	gender	race	religion	age	sentiment	profession	media	nationality	political	multiple social	annotator	sexuality	disability	ethnicity	physical appearance	dialect	human	intersectional	
gender	19.3	8.2	8.8	-8.5	5.8	-8.6	3.0	-6.5	-7.5	-6.5	5.1	1.7	3.8	3.0	1.1	-2.3	-1.7		gender
race	19.3		12.2	12.4	-1.0	6.8	-3.0	4.2	-1.6	-2.6	-2.6	8.6	3.8	4.2	4.6	0.0	-0.8	-0.6	race
religion	8.2	12.2		4.4	-1.4	7.9	-1.0	7.0	-0.9	-0.9	-0.9	4.2	4.2	2.4	5.5	0.7	-0.3	-0.2	religion
age	8.8	12.4	4.4		-1.3	3.9	-1.0	6.1	0.1	-0.9	0.1	4.2	5.3	4.4	5.5	0.7	-0.3	-0.2	age
sentiment	-8.5	-1.0	-1.4	-1.3		0.1	0.1	0.2	-0.7	-0.7	-0.7	-0.7	-0.6	-0.5	-0.4	-0.3	-0.2	0.8	sentiment
profession	5.8	6.8	7.9	3.9	0.1		-0.7	5.4	-0.6	-0.6	-0.6	3.4	4.5	0.6	4.7	-0.2	-0.2	-0.1	profession
media	-8.6	-3.0	-1.0	-1.0	0.1	-0.7		-0.6	-0.6	-0.6	-0.6	-0.5	-0.5	-0.4	-0.3	-0.2	-0.2	-0.1	media
nationality	3.0	4.2	7.0	6.1	0.2	5.4	-0.6		-0.5	-0.5	-0.5	3.5	5.6	3.6	4.7	-0.2	-0.2	-0.1	nationality
political	-6.5	-1.6	-0.9	0.1	-0.7	-0.6	-0.6	-0.5		-0.5	0.5	0.6	0.6	0.7	-0.3	-0.2	-0.1	-0.1	political
multiple social	-7.5	-2.6	-0.9	-0.9	-0.7	-0.6	-0.6	-0.5	-0.5		-0.5	-0.4	-0.4	-0.3	-0.3	-0.2	-0.1	-0.1	multiple social
annotator	-6.5	-2.6	-0.9	0.1	-0.7	-0.6	-0.6	-0.5	0.5	-0.5		-0.4	-0.4	-0.3	-0.3	-0.2	-0.1	-0.1	annotator
sexuality	5.1	8.6	4.2	4.2	-0.7	3.4	-0.5	3.5	0.6	-0.4	-0.4		4.6	0.7	3.8	-0.2	-0.1	-0.1	sexuality
disability	1.7	3.8	4.2	5.3	-0.6	4.5	-0.5	5.6	0.6	-0.4	-0.4	4.6		1.7	4.8	-0.2	-0.1	-0.1	disability
ethnicity	3.8	4.2	2.4	4.4	-0.5	0.6	-0.4	3.6	0.7	-0.3	-0.3	0.7	1.7		0.8	-0.1	-0.1	-0.1	ethnicity
physical appearance	3.0	4.6	5.5	5.5	-0.4	4.7	-0.3	4.7	-0.3	-0.3	-0.3	3.8	4.8	0.8		-0.1	-0.1	-0.1	physical appearance
dialect	1.1	0.0	0.7	0.7	-0.3	-0.2	-0.2	-0.2	-0.2	-0.2	-0.2	-0.2	-0.2	-0.1	-0.1		-0.1	-0.0	dialect
human	-2.3	-0.8	-0.3	-0.3	-0.2	-0.2	-0.2	-0.2	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.0	human
intersectional	-1.7	-0.6	-0.2	-0.2	0.8	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.0	-0.0		intersectional

Table 3: Observed joints: number of papers with combinations of languages (ISO 639), biases

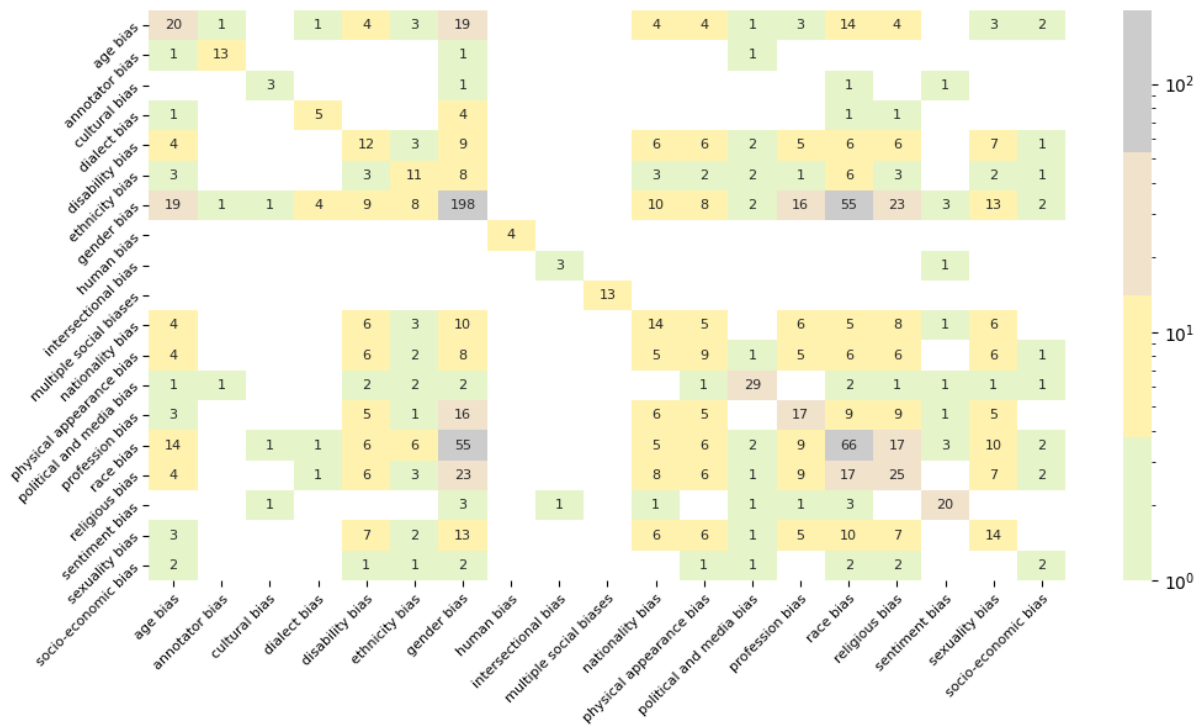


Figure 1: A log-scaled heatmap showing the type and frequency of social biases.

463 itary interests), where researchers should be
 464 aware that their work has broader impact and
 465 can be abused (Hovy and Spruit, 2016).

- 466 • **Availability and Overexposure** We have
 467 found that a small number of papers introduce
 468 new datasets or benchmarks (see 5). Creating

469 and curating new datasets as well as bench-
 470 marks are often a time-consuming, expensive
 471 and long process, where it is oftentimes eas-
 472 ier to utilize existing resources to try out new
 473 methods (Hovy, 2018). Similarly there is the
 474 phenomenon of *topic overexposure*, where
 475 there are waves of seemingly 'popular' re-

en	de	sp	zh	fr	it	ru	tr	pl	nl	ar	pt	da	te	sv	no	ja	hi	sk	ko	he
0.3	2.1	2.9	-4.2	1.6	0.8	-1.2	1.5	-0.7	-1.7	-0.7	0.0	-2.0	-2.2	-1.2	-1.2	-1.2	-2.2	-0.5	0.5	0.5
0.3	2.7	3.7	1.8	1.8	0.8	0.9	0.9	1.9	-0.1	2.9	0.9	-0.1	-0.0	-0.0	-0.0	-0.0	-0.0	1.0	-0.0	1.0
2.1	2.7	8.7	5.0	7.0	4.2	4.5	4.6	1.7	1.7	2.7	2.7	0.7	-0.2	-0.2	-0.2	1.8	-0.2	-0.1	0.9	1.9
2.9	3.7	8.7	3.2	6.2	5.4	3.6	3.7	1.7	-0.3	3.7	3.8	0.8	-0.2	-0.2	-0.2	0.8	-0.2	-0.1	0.9	1.9
-4.2	1.8	5.0	3.2	1.4	0.5	1.7	1.7	0.8	-0.2	2.8	1.8	0.8	-0.1	0.9	-0.1	1.9	-0.1	-0.1	0.9	-0.1
1.6	1.8	7.0	6.2	1.4	6.6	1.7	-0.2	-0.2	0.8	1.8	0.8	0.8	-0.1	-0.1	0.9	-0.1	-0.1	-0.1	-0.1	0.9
0.8	0.8	4.2	5.4	0.5	6.6	1.8	0.8	0.8	0.8	0.8	1.9	0.9	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	0.9
-1.2	0.9	4.5	3.6	1.7	1.7	1.8	1.9	-0.1	-0.1	1.9	0.9	0.9	-0.1	0.9	-0.1	0.9	-0.1	-0.1	-0.0	1.0
1.5	0.9	4.6	3.7	1.7	-0.2	0.8	1.9	1.9	-0.1	-0.1	0.9	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	1.0	-0.0
-0.7	1.9	1.7	1.7	0.8	-0.2	0.8	-0.1	1.9	-0.1	-0.1	1.9	-0.1	-0.0	-0.0	-0.0	-0.0	-0.0	1.0	-0.0	-0.0
-1.7	-0.1	1.7	-0.3	-0.2	0.8	0.8	-0.1	-0.1	-0.1	-0.1	-0.1	-0.1	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
-0.7	2.9	2.7	3.7	2.8	1.8	0.8	1.9	-0.1	-0.1	-0.1	0.9	-0.1	-0.0	-0.0	-0.0	1.0	-0.0	-0.0	-0.0	1.0
0.0	0.9	2.7	3.8	1.8	0.8	1.9	0.9	0.9	1.9	-0.1	0.9	1.0	-0.0	-0.0	-0.0	1.0	-0.0	-0.0	-0.0	-0.0
-2.0	-0.1	0.7	0.8	0.8	0.8	0.9	0.9	-0.1	-0.1	-0.1	1.0	1.0	-0.0	1.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
-2.2	-0.0	-0.2	-0.2	-0.1	-0.1	-0.1	-0.1	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
-1.2	-0.0	-0.2	-0.2	0.9	-0.1	-0.1	0.9	-0.1	-0.0	-0.0	-0.0	-0.0	1.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
-1.2	-0.0	-0.2	-0.2	-0.1	0.9	-0.1	-0.1	-0.1	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
-1.2	-0.0	1.8	0.8	1.9	-0.1	-0.1	0.9	-0.1	-0.0	-0.0	1.0	1.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
-2.2	-0.0	-0.2	-0.2	-0.1	-0.1	-0.1	-0.1	-0.1	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
-0.5	1.0	-0.1	-0.1	-0.1	-0.1	-0.1	-0.0	-0.0	1.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
0.5	-0.0	0.9	0.9	0.9	-0.1	-0.1	-0.0	1.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0
0.5	1.0	1.9	1.9	-0.1	0.9	0.9	1.0	-0.0	-0.0	-0.0	1.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0

Table 4: Residuals between observed joints and predictions based on margins.

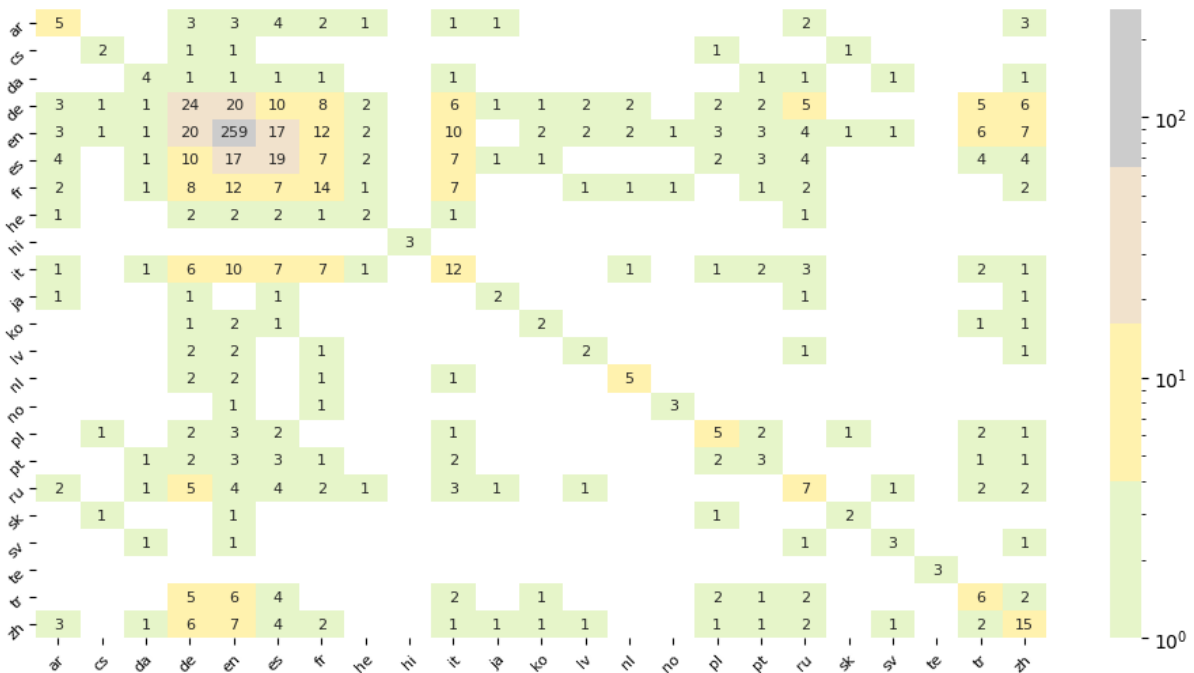


Figure 2: A log-scaled heatmap showing the different languages (ISO-639-1) and their frequency of occurrence.

476 search topics that eventually go out of fashion.
 477 This is based on availability heuristic, if peo-
 478 ple recall a certain event or have knowledge
 479 about certain things then it must be important
 480 (Hovy and Spruit, 2016).

481 **Bias** We have found a limited focus on specific
 482 social biases, where possible causes are rooted in
 483 (i) the data that encodes bias by default (Chandra-
 484 bose et al., 2021), where already available data

485 determines what kind of bias we focus on, (ii) ma-
 486 chine learning breakthroughs in NLP has enabled
 487 ‘streetlamp science’ and we focus on tasks that can
 488 be solved (Hovy, 2018) and (iii) lack of awareness.
 489 This has the consequence that difficult tasks are
 490 not being tackled and bias remains present in NLP
 491 tools. Therefore it is key to raise awareness (Baeza-
 492 Yates, 2018), understand and measure what kind
 493 of bias has influence on NLP models and work to-
 494 wards developing solutions that are equitable. Here,

we (i) showcase challenges in three frequently researched social biases that have been identified through this survey and (ii) point out opportunities for the future with the aim to raise awareness.

- **Gender bias** There is a strong emphasis on a binary understanding of gender (Schnoebelen, 2017; Orgad and Belinkov, 2022) and most task have been reduced to a masculine/feminine dichotomy (Savoldi et al., 2021). Initially, this may be perceived as useful to enable research, however it does not capture the reality of the society and world we live in today. For example, in the USA alone over 1.4 million people identify as transgender (Larson, 2017) and 1.2 million identify as non-binary (Williams Institute, 2021). At the same time, there is not only a common misconception that gender and sex are the same (Larson, 2017), but also that sexuality is somehow indicative of either gender or sex. Sexuality refers to a person’s attraction to a sex or gender, but is not a marker of gender/sex itself (Baum and Westheimer, 2015). Therefore, one can not use sex or gender as a predictor or precursor to assuming a person’s sexuality. Talat et al. (2022b) argue that characteristics like sexuality are usually not observable, which can lead to a reliance on hegemonic stereotypes and unnatural language in bias evaluation benchmarks. This leaves plenty of opportunity to start conversations around developing new datasets, benchmarks and methods that are more inclusive and reflect the world we live in (Savoldi et al., 2021; Orgad and Belinkov, 2022).

- **Race bias** Related work by (Field et al., 2021) provides an excellent overview of the state of the art of race bias research in NLP. In their survey they identify that there are very few datasets and benchmarks and that oftentimes a narrow view of race and racial identity are perpetuated. Additionally, researchers often doesn’t explicitly state if they are focusing on racial bias through downstream tasks such as abusive language detection. Subsequently, currently deployed hate speech or toxicity classifiers mislabel language predominantly used in the African American community as toxic or hate speech when it is not (Dixon et al., 2018; Xia et al., 2020).

Finally, this survey mentions a number of social

biases that have been mentioned such as *religious*, *age* and *disability* with few papers in Figure 1. It is outside of the scope of this work to address each social bias individually, but we would like to point out that there is a lack of relevant benchmarks, datasets and surveys to make substantial progress in these areas and understand the unique challenges each community faces (individually and at an intersectional-level). Most importantly, we would like to emphasize that this type of future work needs to be deeply grounded in interdisciplinary research and led by diverse teams that connect and engage with relevant communities.

Interdisciplinary research The relationship between language and social hierarchies is far more complex than what current techniques can capture. Therefore new methods need to be grounded in relevant literature outside of NLP (Blodgett et al., 2020), because social bias is a complex issue (Sun et al., 2019). Whilst NLP researchers may be committed to using ethical approaches, they may not necessarily have the required ethical and legal knowledge to do so (Santy et al., 2021). This makes it incredibly important to foster collaboration between disciplines to ensure that historical inequalities and biases are taken into consideration when building new algorithms or systems (Caliskan et al., 2017).

Diversity Given the real-life impact of NLP systems and research on people, there is not just a need for diversity in experts working on such systems (Caliskan et al., 2017), but also a need for practitioners and researchers to engage with the affected communities and stakeholders (Blodgett et al., 2020; Fortuna et al., 2021). Therefore, we need to recognize the implicit bias of the people working on different NLP systems and sense-check at different stages how this bias may be reflected in collected data, new benchmarks or models (Savoldi et al., 2021; Hovy and Spruit, 2016). We also need to acknowledge the lack of diversity in teams working on NLP (Schluter, 2018; Savoldi et al., 2021) and work towards more inclusive teams that represent a wide variety of backgrounds and lived experiences (Field et al., 2021). Otherwise, NLP systems continue to represent majorities and we risk the further oppression of already disadvantaged communities (Talat et al., 2022b; Schnoebelen, 2017).

495
496
497
498

499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528

529
530
531
532
533
534
535
536
537
538
539
540
541
542
543

544

545
546
547
548
549
550
551
552
553
554
555
556
557

558
559
560
561
562
563
564
565
566
567
568
569
570
571
572

573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592

593 Limitations and Ethics Statement

594 In this paper, we surveyed a collection of papers
595 and identified continued challenges in social bias
596 research. We have created this collection based on
597 a keyword search and outlined how this may not
598 fully reflect all literature on social bias existing in
599 the ACL anthology or other venues. We only used
600 open-access papers in this collection and no human
601 participants were involved in this work. Tradition-
602 ally, social biases have been investigated in fields
603 such as social sciences, law, or psychology which
604 we have not discussed here. Furthermore, we do
605 not give an analysis of algorithmic or dataset biases
606 (e.g., machine learning, data mining or otherwise)
607 or provided an in-depth review of technical con-
608 tributions in computational social biases. We are
609 also limited by the resource of the reviewed papers,
610 where substantial contributions to the field have
611 been made outside of ACL venues. Finally, we
612 would like to point out that opportunities and rec-
613 ommendations for future bias research as proposed
614 in section 6 should be considered from a euro- and/
615 or anglo-centric perspective. There may be a vari-
616 ation depending on the social context, country or
617 culture that works on a specific bias problem.

618 Acknowledgements

619 References

620 Samyak Agrawal, Kshitij Gupta, Devansh Gautam, and
621 Radhika Mamidi. 2022. Towards detecting political
622 bias in hindi news articles. In *Proceedings of the 60th*
623 *Annual Meeting of the Association for Computational*
624 *Linguistics: Student Research Workshop*, pages 239–
625 244.

626 Carlos Aguirre, Keith Harrigan, and Mark Dredze.
627 2021. Gender and racial fairness in depression re-
628 search using social media. In *Proceedings of the 16th*
629 *Conference of the European Chapter of the Associ-*
630 *ation for Computational Linguistics: Main Volume*,
631 pages 2932–2949.

632 Jaimeen Ahn, Hwaran Lee, Jinhwa Kim, and Alice Oh.
633 2022. Why knowledge distillation amplifies gender
634 bias and how to mitigate from the perspective of dis-
635 tilbert. In *Proceedings of the 4th Workshop on Gen-*
636 *der Bias in Natural Language Processing (GeBNLP)*,
637 pages 266–272.

638 Jaimeen Ahn and Alice Oh. 2021. Mitigating language-
639 dependent ethnic bias in bert. In *Proceedings of the*
640 *2021 Conference on Empirical Methods in Natural*
641 *Language Processing*, pages 533–549.

642 Dmitrii Aksenov, Peter Bourgonje, Karolina Zaczyn-
643 ska, Malte Ostendorff, Julian Moreno Schneider, and

Georg Rehm. 2021. Fine-grained classification of
political bias in german news: A data set and initial
experiments. In *Proceedings of the 5th Workshop*
on Online Abuse and Harms (WOAH 2021), pages
121–131.

Afra Feyza Akyürek, Muhammed Yusuf Kocyigit, Se-
jin Paik, and Derry Wijaya. 2022a. Challenges in
measuring bias via open-ended language generation.
arXiv preprint arXiv:2205.11601.

Afra Feyza Akyürek, Muhammed Yusuf Kocyigit, Sejin
Paik, and Derry Tanti Wijaya. 2022b. [Challenges in
measuring bias via open-ended language generation](#).
In *Proceedings of the 4th Workshop on Gender Bias*
in Natural Language Processing (GeBNLP), pages
76–76, Seattle, Washington. Association for Compu-
tational Linguistics.

Hala Al Kuwatly, Maximilian Wich, and Georg Groh.
2020. Identifying and measuring annotator bias
based on annotators’ demographic characteristics. In
Proceedings of the Fourth Workshop on Online Abuse
and Harms, pages 184–190.

Laura Alonso Alemany, Luciana Benotti, Hernán
Maina, Lucía González, Lautaro Martínez, Beatriz
Busaniche, Alexia Halvorsen, Amanda Rojo, and
Mariela Rajngewerc. 2023. Bias assessment for ex-
perts in discrimination, not in computer science. In
Proceedings of the First Workshop on Cross-Cultural
Considerations in NLP (C3NLP), pages 91–106.

Sarah Alnegheimish, Alicia Guo, and Yi Sun. 2022.
Using natural sentence prompts for understanding bi-
ases in language models. In *Proceedings of the 2022*
Conference of the North American Chapter of the
Association for Computational Linguistics: Human
Language Technologies, pages 2824–2830.

Saied Alshahrani, Esmá Wali, Abdullah R Alshamsan,
Yan Chen, and Jeanna Matthews. 2022. Roadblocks
in gender bias measurement for diachronic corpora.
In *Proceedings of the 3rd Workshop on Computa-*
tional Approaches to Historical Language Change,
pages 140–148.

Karl Ameriks and Desmond M Clarke. 2000. *Aristotle:*
Nicomachean Ethics. Cambridge University Press.

Haozhe An, Zongxia Li, Jieyu Zhao, and Rachel
Rudinger. 2022a. Sodapop: Open-ended discovery
of social biases in social commonsense reasoning
models. *arXiv preprint arXiv:2210.07269*.

Haozhe An, Xiaojiang Liu, and Donald Zhang. 2022b.
Learning bias-reduced word embeddings using dic-
tionary definitions. In *Findings of the Association for*
Computational Linguistics: ACL 2022, pages 1139–
1152.

Talita Anthonio and Lennart Kloppenburg. 2019. Team
kermit-the-frog at semeval-2019 task 4: Bias detec-
tion through sentiment analysis and simple linguistic
features. In *Proceedings of the 13th International*
Workshop on Semantic Evaluation, pages 1016–1020.

700	Maria Antoniak and David Mimno. 2021. Bad seeds:	Francesco Barbieri, Francesco Ronzano, and Horacio	755
701	Evaluating lexical methods for bias measurement.	Saggion. 2015. How topic biases your results? a case	756
702	In <i>Proceedings of the 59th Annual Meeting of the</i>	study of sentiment analysis and irony detection in ital-	757
703	<i>Association for Computational Linguistics and the</i>	ian. In <i>Proceedings of the International Conference</i>	758
704	<i>11th International Joint Conference on Natural Lan-</i>	<i>Recent Advances in Natural Language Processing</i> ,	759
705	<i>guage Processing (Volume 1: Long Papers)</i> , pages	pages 41–47.	760
706	1889–1904.		
707	Giuseppe Attanasio, Debora Nozza, Dirk Hovy, and	Soumya Barikeri, Anne Lauscher, Ivan Vulić, and Goran	761
708	Elena Baralis. 2022. Entropy-based attention regu-	Glavaš. 2021. Reddittbias: A real-world resource for	762
709	larization frees unintended bias mitigation from lists.	bias evaluation and debiasing of conversational lan-	763
710	In <i>Findings of the Association for Computational</i>	guage models. In <i>Proceedings of the 59th Annual</i>	764
711	<i>Linguistics: ACL 2022</i> , pages 1105–1119.	<i>Meeting of the Association for Computational Lin-</i>	765
712		<i>guistics and the 11th International Joint Conference</i>	766
713	Hossein Azarpanah and Mohsen Farhadloo. 2021. Mea-	<i>on Natural Language Processing (Volume 1: Long</i>	767
714	suring biases of word embeddings: What similar-	<i>Papers)</i> , pages 1941–1955.	768
715	ity measures and descriptive statistics to use? In	Marion Bartl, Malvina Nissim, and Albert Gatt. 2020.	769
716	<i>Proceedings of the First Workshop on Trustworthy</i>	Unmasking contextual stereotypes: Measuring and	770
717	<i>Natural Language Processing</i> , pages 8–14.	mitigating bert’s gender bias. In <i>Proceedings of the</i>	771
718		<i>Second Workshop on Gender Bias in Natural Lan-</i>	772
719	Ricardo Baeza-Yates. 2018. Bias on the Web. <i>Communi-</i>	<i>guage Processing</i> , pages 1–16.	773
720	cations of the ACM, 61(6):54–61.	Christine Basta, Marta R Costa-jussà, and Noe Casas.	774
721	Sunyam Bagga and Andrew Piper. 2020. Measuring the	2019. Evaluating the underlying gender bias in con-	775
722	effects of bias in training data for literary classifica-	textualized word embeddings. In <i>Proceedings of the</i>	776
723	tion. In <i>Proceedings of the The 4th Joint SIGHUM</i>	<i>First Workshop on Gender Bias in Natural Language</i>	777
724	<i>Workshop on Computational Linguistics for Cultural</i>	<i>Processing</i> , pages 33–39.	778
725	<i>Heritage, Social Sciences, Humanities and Literature</i> ,	Lisa Bauer, Hanna Tischer, and Mohit Bansal. 2023. So-	779
726	pages 74–84.	cial commonsense for explanation and cultural bias	780
727	Esma Balkir, Svetlana Kiritchenko, Isar Nejadgholi,	discovery. In <i>Proceedings of the 17th Conference of</i>	781
728	and Kathleen Fraser. 2022. Challenges in apply-	<i>the European Chapter of the Association for Computa-</i>	782
729	ing explainability methods to improve the fairness	<i>tional Linguistics</i> , pages 3727–3742.	783
730	of NLP models. In <i>Proceedings of the 2nd Work-</i>	Joel Baum and Kim Westheimer. 2015. Sex? Sexual	784
731	<i>shop on Trustworthy Natural Language Processing</i>	Orientation? Gender Identity? Gender Expression?	785
732	<i>(TrustNLP 2022)</i> , pages 80–92, Seattle, U.S.A. Asso-	<i>Teaching Tolerance</i> , 50(34–38).	786
733	ciation for Computational Linguistics.	Emily M Bender. 2011. On achieving and evaluating	787
734	Ramy Baly, Giovanni Da San Martino, James Glass,	language-independence in NLP. <i>Linguistic Issues in</i>	788
735	and Preslav Nakov. 2020a. We can detect your bias:	<i>Language Technology</i> , 6.	789
736	Predicting the political ideology of news articles. In	Emily M Bender and Batya Friedman. 2018. Data	790
737	<i>Proceedings of the 2020 Conference on Empirical</i>	statements for natural language processing: Toward	791
738	<i>Methods in Natural Language Processing (EMNLP)</i> ,	mitigating system bias and enabling better science.	792
739	pages 4982–4991.	<i>Transactions of the Association for Computational</i>	793
740	Ramy Baly, Georgi Karadzhov, Dimitar Alexandrov,	<i>Linguistics</i> , 6:587–604.	794
741	James Glass, and Preslav Nakov. 2018. Predict-	Emily M Bender, Dirk Hovy, and Alexandra Schofield.	795
742	ing factuality of reporting and bias of news media	2020. Integrating ethics into the nlp curriculum. In	796
743	sources. In <i>Proceedings of the 2018 Conference on</i>	<i>Proceedings of the 58th Annual Meeting of the As-</i>	797
744	<i>Empirical Methods in Natural Language Processing</i> ,	<i>sociation for Computational Linguistics: Tutorial</i>	798
745	pages 3528–3539.	<i>Abstracts</i> , pages 6–9.	799
746	Ramy Baly, Georgi Karadzhov, Jisun An, Haewoon	Luciana Benotti and Patrick Blackburn. 2022. Ethics	800
747	Kwak, Yoan Dinkov, Ahmed Ali, James Glass, and	consideration sections in natural language processing	801
748	Preslav Nakov. 2020b. What was written vs. who	papers. In <i>Proceedings of the 2022 Conference on</i>	802
749	read it: News media profiling using text analysis	<i>Empirical Methods in Natural Language Processing</i> ,	803
750	and social media context. In <i>Proceedings of the 58th</i>	pages 4509–4516.	804
751	<i>Annual Meeting of the Association for Computational</i>	Luciana Benotti, Karën Fort, Min-Yen Kan, and Yulia	805
752	<i>Linguistics</i> , pages 3364–3374.	Tsvetkov. 2023. Understanding ethics in nlp author-	806
753	Hritik Bansal, Da Yin, Masoud Monajatipoor, and Kai-	ing and reviewing. In <i>Proceedings of the 17th Con-</i>	807
754	Wei Chang. 2022. How well can text-to-image gen-	<i>ference of the European Chapter of the Association</i>	808
	erative models understand ethical natural language	<i>for Computational Linguistics: Tutorial Abstracts</i> ,	809
	interventions? <i>arXiv preprint arXiv:2210.15230</i> .	pages 19–24.	810

811	Adrian Benton, Glen Coppersmith, and Mark Dredze.	Nadine Braun, Martijn Goudbeek, and Emiel Kra-	866
812	2017. Ethical research protocols for social media	mer. 2016. The multilingual affective soccer corpus	867
813	health research. In <i>Proceedings of the first ACL work-</i>	(masc): Compiling a biased parallel corpus on soc-	868
814	<i>shop on ethics in natural language processing</i> , pages	cer reportage in english, german and dutch. In <i>Pro-</i>	869
815	94–102.	<i>ceedings of the 9th International Natural Language</i>	870
		<i>Generation Conference</i> , pages 74–78.	871
816	Amanda Bertsch, Ashley Oh, Sanika Natu, Swetha	Luke Breitfeller, Emily Ahn, David Jurgens, and Yulia	872
817	Gangu, Alan W Black, and Emma Strubell. 2022.	Tsvetkov. 2019. Finding microaggressions in the	873
818	Evaluating gender bias transfer from film data. In	wild: A case for locating elusive phenomena in social	874
819	<i>Proceedings of the 4th Workshop on Gender Bias</i>	media posts. In <i>Proceedings of the 2019 conference</i>	875
820	<i>in Natural Language Processing (GeBNLP)</i> , pages	<i>on empirical methods in natural language processing</i>	876
821	235–243.	<i>and the 9th international joint conference on natural</i>	877
822	Advait Bhat, Saaket Agashe, and Anirudha Joshi. 2021.	<i>language processing (EMNLP-IJCNLP)</i> , pages 1664–	878
823	How do people interact with biased text prediction	1674.	879
824	models while writing? In <i>Proceedings of the First</i>		
825	<i>Workshop on Bridging Human–Computer Interaction</i>	Pere-Lluís Huguet Cabot, David Abadi, Agneta Fischer,	880
826	<i>and Natural Language Processing</i> , pages 116–121.	and Ekaterina Shutova. 2021. Us vs. them: A dataset	881
		of populist attitudes, news bias and emotions. In	882
827	Su Lin Blodgett, Solon Barocas, Hal Daumé III, and	<i>Proceedings of the 16th Conference of the European</i>	883
828	Hanna Wallach. 2020. Language (technology) is	<i>Chapter of the Association for Computational Lin-</i>	884
829	power: A critical survey of “bias” in NLP . In <i>Pro-</i>	<i>guistics: Main Volume</i> , pages 1921–1945.	885
830	<i>ceedings of the 58th Annual Meeting of the Associa-</i>		
831	<i>tion for Computational Linguistics</i> , pages 5454–5476.	Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan.	886
832	Association for Computational Linguistics.	2017. Semantics derived automatically from lan-	887
		guage corpora contain human-like biases. <i>Science</i> ,	888
833	Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu,	356(6334):183–186.	889
834	Robert Sim, and Hanna Wallach. 2021. Stereotyping	António Câmara, Nina Taneja, Tamjeed Azad, Emily	890
835	norwegian salmon: An inventory of pitfalls in fair-	Allaway, and Richard Zemel. 2022. Mapping the	891
836	ness benchmark datasets. In <i>Proceedings of the 59th</i>	multilingual margins: Intersectional biases of senti-	892
837	<i>Annual Meeting of the Association for Computational</i>	ment analysis systems in english, spanish, and ara-	893
838	<i>Linguistics and the 11th International Joint Confer-</i>	bic. In <i>Proceedings of the Second Workshop on Lan-</i>	894
839	<i>ence on Natural Language Processing (Volume 1:</i>	<i>guage Technology for Equality, Diversity and Inclu-</i>	895
840	<i>Long Papers)</i> , pages 1004–1015.	<i>sion</i> , pages 90–106.	896
841	Tolga Bolukbasi, Kai-Wei Chang, James Y Zou,	Yang Cao, Yada Pruksachatkun, Kai-Wei Chang, Rahul	897
842	Venkatesh Saligrama, and Adam T Kalai. 2016. Man	Gupta, Varun Kumar, Jwala Dhamala, and Aram Gal-	898
843	is to computer programmer as woman is to home-	styan. 2022a. On the intrinsic and extrinsic fairness	899
844	maker? Debiasing word embeddings. <i>Advances in</i>	evaluation metrics for contextualized language repre-	900
845	<i>neural information processing systems</i> , 29.	sentations. In <i>Proceedings of the 60th Annual Meet-</i>	901
		<i>ing of the Association for Computational Linguistics</i>	902
846	Conrad Borchers, Dalia Gala, Benjamin Gilbert, Eduard	<i>(Volume 2: Short Papers)</i> , pages 561–570.	903
847	Oravkin, Wilfried Bounsi, Yuki M Asano, and Han-	Yang Cao, Anna Sotnikova, Hal Daumé III, Rachel	904
848	nah Kirk. 2022. Looking for a handsome carpenter!	Rudinger, and Linda Zou. 2022b. Theory-grounded	905
849	debiasing gpt-3 job advertisements. In <i>Proceedings</i>	measurement of us social stereotypes in english lan-	906
850	<i>of the 4th Workshop on Gender Bias in Natural Lan-</i>	guage models. In <i>Proceedings of the 2022 Confer-</i>	907
851	<i>guage Processing (GeBNLP)</i> , pages 212–224.	<i>ence of the North American Chapter of the Associ-</i>	908
		<i>ation for Computational Linguistics: Human Lan-</i>	909
852	Shikha Bordia and Samuel Bowman. 2019. Identifying	<i>guage Technologies</i> , pages 1276–1295.	910
853	and reducing gender bias in word-level language	Yang Trista Cao and Hal Daumé III. 2021. Toward	911
854	models. In <i>Proceedings of the 2019 Conference of</i>	gender-inclusive coreference resolution: An analysis	912
855	<i>of the North American Chapter of the Association for</i>	of gender and bias throughout the machine learn-	913
856	<i>Computational Linguistics: Student Research Work-</i>	ing lifecycle. <i>Computational Linguistics</i> , 47(3):615–	914
857	<i>shop</i> , pages 7–15.	661.	915
858	Tom Bourgeade, Alessandra Teresa Cignarella, Si-	Ilias Chalkidis, Tommaso Pasini, Sheng Zhang, Letizia	916
859	mona Frenda, Mario Laurent, Wolfgang Schmeisser-	Tomada, Sebastian Schwemer, and Anders Søgaard.	917
860	Nieto, Farah Benamara, Cristina Bosco, Véronique	2022. Fairlex: A multilingual benchmark for evaluat-	918
861	Moriceau, Viviana Patti, and Mariona Taulé. 2023.	ing fairness in legal text processing. In <i>Proceedings</i>	919
862	A multilingual dataset of racial stereotypes in social	<i>of the 60th Annual Meeting of the Association for</i>	920
863	media conversational threads. In <i>Findings of the As-</i>	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	921
864	<i>sociation for Computational Linguistics: EACL 2023</i> ,	pages 4389–4406.	922
865	pages 674–684.		

923	Kaytlin Chaloner and Alfredo Maldonado. 2019. Measuring gender bias in word embeddings across domains and discovering new gender bias word categories. In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 25–32.	978
924		979
925		980
926		981
927		
928		
929	Aravindan Chandrabose, Bharathi Raja Chakravarthi, et al. 2021. An overview of fairness in data-illuminating the bias in data pipeline. In <i>Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion</i> , pages 34–45.	982
930		983
931		984
932		985
933		986
934	Kai-Wei Chang, Vinodkumar Prabhakaran, and Vicente Ordonez. 2019. Bias and fairness in natural language processing. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): Tutorial Abstracts</i> .	987
935		988
936		989
937		990
938		991
939		992
940		
941	John Chen, Ian Berlot-Attwell, Xindi Wang, Safwan Hossain, and Frank Rudzicz. 2020a. Exploring text specific and blackbox fairness algorithms in multimodal clinical nlp. In <i>Proceedings of the 3rd Clinical Natural Language Processing Workshop</i> , pages 301–312.	993
942		994
943		995
944		996
945		997
946		998
947	Wei-Fan Chen, Khalid Al Khatib, Benno Stein, and Henning Wachsmuth. 2020b. Detecting media bias in news articles using gaussian bias distributions. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4290–4300.	999
948		1000
949		1001
950		1002
951		
952	Wei-Fan Chen, Henning Wachsmuth, Khalid Al Khatib, and Benno Stein. 2018. Learning to flip the bias of news headlines. In <i>Proceedings of the 11th International conference on natural language generation</i> , pages 79–88.	1003
953		1004
954		1005
955		1006
956		1007
957	Xiuying Chen, Mingzhe Li, Rui Yan, Xin Gao, and Xiangliang Zhang. 2022. Unsupervised mitigating gender bias by character components: A case study of chinese word embedding. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 121–128.	1008
958		1009
959		1010
960		1011
961		1012
962		
963	Lu Cheng, Ahmadreza Mosallanezhad, Yasin Silva, Deborah Hall, and Huan Liu. 2021. Mitigating bias in session-based cyberbullying detection: A non-compromising approach. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2158–2168.	1013
964		1014
965		1015
966		1016
967		1017
968		
969		
970		
971	Won Ik Cho, Ji Won Kim, Seok Min Kim, and Nam Soo Kim. 2019. On measuring gender bias in translation of gender-neutral pronouns. In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 173–181.	1018
972		1019
973		1020
974		1021
975		1022
976	Jonathan Gabel Christiansen, Mathias Gammelgaard, and Anders Sjøgaard. 2021. The effect of round-trip translation on fairness in sentiment analysis. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4423–4428.	1023
977		1024
		1025
		1026
	Yung-Sung Chuang, Mingye Gao, Hongyin Luo, James Glass, Hung-Yi Lee, Yun-Nung Chen, and Shang-Wen Li. 2021. Mitigating biases in toxic language detection through invariant rationalization. In <i>Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021)</i> , pages 114–120.	1027
		1028
		1029
		1030
		1031
	Chloe Ciora, Nur Iren, and Malihe Alikhani. 2021. Examining covert gender bias: A case study in turkish and english machine translation models. In <i>Proceedings of the 14th International Conference on Natural Language Generation</i> , pages 55–63.	
	Trevor Cohn and Lucia Specia. 2013. Modelling annotator bias with multi-task gaussian processes: An application to machine translation quality estimation. In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 32–42.	
	Kimberlé Williams Crenshaw. 2013. Mapping the margins: Intersectionality, identity politics, and violence against women of color. In <i>The public nature of private violence</i> , pages 93–118. Routledge.	
	Paula Czarowska, Yogarshi Vyas, and Kashif Shah. 2021. Quantifying social biases in NLP: A generalization and empirical comparison of extrinsic fairness metrics. <i>Transactions of the Association for Computational Linguistics</i> , 9:1249–1267.	
	Joke Daems and Janiça Hackenbuchner. 2022. Debi-asbyus: Raising awareness and creating a database of mt bias. In <i>Proceedings of the 23rd Annual Conference of the European Association for Machine Translation</i> , pages 287–288.	
	Dipto Das, Shion Guha, and Bryan Semaan. 2023. Toward cultural bias evaluation datasets: The case of bengali gender, religious, and national identity. In <i>Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)</i> , pages 68–83.	
	Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. 2019. Racial bias in hate speech and abusive language detection datasets. In <i>Proceedings of the Third Workshop on Abusive Language Online</i> , pages 25–35.	
	Hillary Dawkins. 2021a. Marked attribute bias in natural language inference. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4214–4226.	
	Hillary Dawkins. 2021b. Second order winobias (sowinobias) test set for latent gender bias detection in coreference resolution. In <i>Proceedings of the 3rd Workshop on Gender Bias in Natural Language Processing</i> , pages 103–111.	

1032	Erenay Dayanik and Sebastian Padó. 2020. Masking actor information leads to fairer political claims detection. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4385–4391.		
1033			
1034			
1035			
1036			
1037	FMG De Jong, Bente Maegaard, Koenraad De Smedt, Darja Fišer, Dieter Van Uytvanck, et al. 2018. Clarin: towards fair and responsible data science using language resources. In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , pages 3259–3264.		
1038			
1039			
1040			
1041			
1042			
1043	Daniel de Vassimon Manela, David Errington, Thomas Fisher, Boris van Breugel, and Pasquale Minervini. 2021. Stereotype and skew: Quantifying gender bias in pre-trained and fine-tuned language models. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2232–2242.		
1044			
1045			
1046			
1047			
1048			
1049			
1050	Pieter Delobelle, Ewoenam Tokpo, Toon Calders, and Bettina Berendt. 2022. Measuring fairness with biased rulers: A comparative study on bias metrics for pre-trained language models. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1693–1706.		
1051			
1052			
1053			
1054			
1055			
1056			
1057	Sunipa Dev, Tao Li, Jeff M Phillips, and Vivek Sriku-mar. 2021a. Oscar: Orthogonal subspace correction and rectification of biases in word embeddings. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5034–5050.		
1058			
1059			
1060			
1061			
1062			
1063	Sunipa Dev, Emily Sheng, Jieyu Zhao, Aubrie Amstutz, Jiao Sun, Yu Hou, Mattie Sansverino, Jiin Kim, Akihiro Nishi, Nanyun Peng, and Kai-Wei Chang. 2021b. What do bias measures measure?		
1064			
1065			
1066			
1067	Hannah Devinney, Jenny Björklund, and Henrik Björklund. 2020. Semi-supervised topic modeling for gender bias discovery in english and swedish. In <i>Proceedings of the Second Workshop on Gender Bias in Natural Language Processing</i> , pages 79–92.		
1068			
1069			
1070			
1071			
1072	Emily Dinan, Angela Fan, Adina Williams, Jack Urbanek, Douwe Kiela, and Jason Weston. 2020a. Queens are powerful too: Mitigating gender bias in dialogue generation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 8173–8188.		
1073			
1074			
1075			
1076			
1077			
1078	Emily Dinan, Angela Fan, Ledell Wu, Jason Weston, Douwe Kiela, and Adina Williams. 2020b. Multi-dimensional gender bias classification. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 314–331.		
1079			
1080			
1081			
1082			
1083			
1084	Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In <i>Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society</i> , pages 67–73.		
1085			
1086			
1087			
1088			
	Jad Doughman, Wael Khreich, Maya El Gharib, Maha Wiss, and Zahraa Berjawi. 2021. Gender bias in text: Origin, taxonomy, and implications. In <i>Proceedings of the 3rd Workshop on Gender Bias in Natural Language Processing</i> , pages 34–44.	1089	1090
		1091	1092
		1093	1094
	Jo Drugan and Bogdan Babych. 2010. Shared resources, shared values? ethical implications of sharing translation resources. In <i>Proceedings of the Second Joint EM+/CNGL Workshop: Bringing MT to the User: Research on Integrating MT in the Translation Industry</i> , pages 3–10.	1094	1095
		1096	1097
		1098	1099
		1100	1101
	Yupei Du, Qixiang Fang, and Dong Nguyen. 2021. Assessing the reliability of word embedding gender bias measures. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10012–10034.	1102	1103
		1104	1105
	Yupei Du, Qi Zheng, Yuanbin Wu, Man Lan, Yan Yang, and Meirong Ma. 2022. Understanding gender bias in knowledge base embeddings. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1381–1395.	1106	1107
		1108	1109
		1110	1111
	Penelope Eckert and John R Rickford. 2001. <i>Style and sociolinguistic variation</i> . Cambridge University Press.	1112	1113
		1114	1115
	Fatma Elsafoury, Steven R Wilson, and Naeem Ramzan. 2022. A comparative study on word embeddings and social nlp tasks. In <i>Proceedings of the Tenth International Workshop on Natural Language Processing for Social Media</i> , pages 55–64.	1116	1117
		1118	1119
	Carla Parra Escartin, Wessel Reijers, Teresa Lynn, Joss Moorkens, Andy Way, and Chao-Hong Liu. 2017. Ethical considerations in nlp shared tasks. <i>EACL 2017</i> , page 66.	1120	1121
		1122	1123
	Kawin Ethayarajh. 2020. Is your classifier actually biased? measuring fairness under uncertainty with bernstein bounds. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2914–2919.	1124	1125
		1126	1127
	Agnieszka Falenska and Özlem Çetinoğlu. 2021. Assessing gender bias in Wikipedia: Inequalities in article titles. In <i>Proceedings of the 3rd Workshop on Gender Bias in Natural Language Processing</i> , pages 75–85, Online. Association for Computational Linguistics.	1128	1129
		1130	1131
		1132	1133
	Angela Fan and Claire Gardent. 2022. Generating biographies on wikipedia: The impact of gender bias on the retrieval-based generation of women biographies. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8561–8576.	1134	1135
		1136	1137
		1138	1139
	Lisa Fan, Marshall White, Eva Sharma, Ruisi Su, Prafulla Kumar Choubey, Ruihong Huang, and Lu Wang. 2019. In plain sight: Media bias through the lens of	1140	1141
		1142	

1143	factual reporting. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6343–6349.		
1144			
1145			
1146			
1147			
1148	Michael Färber, Agon Qurdina, and Lule Ahmedi. 2019. Team peter brinkmann at semeval-2019 task 4: Detecting biased news articles using convolutional neural networks. In <i>Proceedings of the 13th International Workshop on Semantic Evaluation</i> , pages 1032–1036.		
1149			
1150			
1151			
1152			
1153			
1154	Anjalie Field, Su Lin Blodgett, Zeerak Waseem, and Yulia Tsvetkov. 2021. A survey of race, racism, and anti-racism in NLP. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1905–1925.		
1155			
1156			
1157			
1158			
1159			
1160			
1161	Anjalie Field and Yulia Tsvetkov. 2020. Unsupervised discovery of implicit gender bias. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 596–608.		
1162			
1163			
1164			
1165			
1166	Joseph Fisher, Arpit Mittal, Dave Palfrey, and Christos Christodoulopoulos. 2020. Debiasing knowledge graph embeddings. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7332–7345.		
1167			
1168			
1169			
1170			
1171	Lucie Flekova, Jordan Carpenter, Salvatore Giorgi, Lyle Ungar, and Daniel Preoțiuc-Pietro. 2016. Analyzing biases in human perception of user age and gender from text. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 843–854.		
1172			
1173			
1174			
1175			
1176			
1177	Joel Escudé Font and Marta R Costa-jussà. 2019. Equalizing gender bias in neural machine translation with word embeddings techniques. In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 147–154.		
1178			
1179			
1180			
1181			
1182	Karën Fort and Alain Couillault. 2016. Yes, we care! results of the ethics and natural language processing surveys. In <i>international Language Resources and Evaluation Conference (LREC) 2016</i> .		
1183			
1184			
1185			
1186	Paula Fortuna, Laura Pérez-Mayos, Leo Wanner, et al. 2021. Cartography of natural language processing for social good (nlp4sg): Searching for definitions, statistics and white spots. In <i>Proceedings of the 1st Workshop on NLP for Positive Impact</i> , pages 19–26.		
1187			
1188			
1189			
1190			
1191	Scott Friedman, Sonja Schmer-Galunder, Anthony Chen, and Jeffrey Rye. 2019. Relating word embedding gender biases to gender gaps: A cross-cultural analysis. In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 18–24.		
1192			
1193			
1194			
1195			
1196			
	Annemarie Friedrich and Torsten Zesch. 2021. A crash course on ethics for natural language processing. In <i>Proceedings of the Fifth Workshop on Teaching NLP</i> , pages 49–51.	1197	1198
	Niklas Friedrich, Anne Lauscher, Simone Paolo Ponzetto, and Goran Glavaš. 2021. Debie: A platform for implicit and explicit debiasing of word embedding spaces. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations</i> , pages 91–98.	1201	1202
	Zee Fryer, Vera Axelrod, Ben Packer, Alex Beutel, Jilin Chen, and Kellie Webster. 2022. Flexible text generation for counterfactual fairness probing. In <i>Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH)</i> , pages 209–229.	1208	1209
	Yacine Gaci, Boualem Benattallah, Fabio Casati, and Khalid Benabdeslem. 2022. Debiasing pretrained text encoders by paying attention to paying attention. In <i>2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9582–9602. Association for Computational Linguistics.	1213	1214
	Marco Gaido, Beatrice Savoldi, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. How to split: the effect of word segmentation on gender bias in speech translation. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3576–3589.	1219	1220
	Dhruvil Gala, Mohammad Omar Khursheed, Hannah Lerner, Brendan O’Connor, and Mohit Iyyer. 2020. Analyzing gender bias within narrative tropes. In <i>Proceedings of the Fourth Workshop on Natural Language Processing and Computational Social Science</i> , pages 212–217.	1225	1226
	Rama Rohit Reddy Gangula, Suma Reddy Duggenpudi, and Radhika Mamidi. 2019. Detecting political bias in news articles using headline attention. In <i>Proceedings of the 2019 ACL workshop BlackboxNLP: analyzing and interpreting neural networks for NLP</i> , pages 77–84.	1231	1232
	Aparna Garimella, Akhash Amarnath, Kiran Kumar, Akash Pramod Yalla, N Anandhavelu, Niyati Chhaya, and Balaji Vasan Srinivasan. 2021. He is very intelligent, she is very beautiful? on mitigating social biases in language modelling and generation. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4534–4545.	1237	1238
	Aparna Garimella, Carmen Banea, Dirk Hovy, and Rada Mihalcea. 2019. Women’s syntactic resilience and men’s grammatical luck: Gender-bias in part-of-speech tagging and dependency parsing. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3493–3498.	1244	1245
	Ismael Garrido-Muñoz, Arturo Montejo-Ráez, Fernando Martínez-Santiago, and L Alfonso Ureña-López. 2021. A survey on bias in deep NLP. <i>Applied Sciences</i> , 11(7):3184.	1250	1251

1254	Andrew Gaut, Tony Sun, Shirlyn Tang, Yuxin Huang,	detection models. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> ,	1310
1255	Jing Qian, Mai ElSherief, Jieyu Zhao, Diba Mirza,	pages 4873–4885.	1311
1256	Elizabeth Belding, Kai-Wei Chang, et al. 2020. To-		1312
1257	wards understanding gender bias in relation extrac-		
1258	tion. In <i>Proceedings of the 58th Annual Meeting of</i>	Yue Guo, Yi Yang, and Ahmed Abbasi. 2022. Auto-	1313
1259	<i>the Association for Computational Linguistics</i> , pages	debias: Debiasing masked language models with	1314
1260	2943–2953.	automated biased prompts. In <i>Proceedings of the</i>	1315
		<i>60th Annual Meeting of the Association for Compu-</i>	1316
1261	Mor Geva, Yoav Goldberg, and Jonathan Berant. 2019.	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	1317
1262	Are we modeling the task or the annotator? an inves-	1012–1023.	1318
1263	tigation of annotator bias in natural language under-		
1264	standing datasets. In <i>Proceedings of the 2019 Confer-</i>	Saurabh Gupta, Hong Huy Nguyen, Junichi Yamagishi,	1319
1265	<i>ence on Empirical Methods in Natural Language Pro-</i>	and Isao Echizen. 2020. Viable threat on news read-	1320
1266	<i>cessing and the 9th International Joint Conference</i>	ing: Generating biased news using natural language	1321
1267	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	models. In <i>Proceedings of the Fourth Workshop on</i>	1322
1268	pages 1161–1166.	<i>Natural Language Processing and Computational</i>	1323
		<i>Social Science</i> , pages 55–65.	1324
1269	Sayan Ghosh, Dylan Baker, David Jurgens, and Vin-		
1270	odkumar Prabhakaran. 2021. Detecting cross-	Umang Gupta, Jwala Dhamala, Varun Kumar, Apurv	1325
1271	geographic biases in toxicity modeling on social me-	Verma, Yada Pruksachatkun, Satyapriya Krishna,	1326
1272	dia. In <i>Proceedings of the Seventh Workshop on</i>	Rahul Gupta, Kai-Wei Chang, Greg Ver Steeg, and	1327
1273	<i>Noisy User-generated Text (W-NUT 2021)</i> , pages 313–	Aram Galstyan. 2022. Mitigating gender bias in dis-	1328
1274	328.	tilled language models via counterfactual role rever-	1329
		sal. In <i>Findings of the Association for Computational</i>	1330
1275	Noa Baker Gillis. 2021. Sexism in the judiciary: The	<i>Linguistics: ACL 2022</i> , pages 658–678.	1331
1276	importance of bias definition in nlp and in our courts.		
1277	In <i>Proceedings of the 3rd Workshop on Gender Bias</i>	Gerhard Hagerer, David Szabo, Andreas Koch, Maria	1332
1278	<i>in Natural Language Processing</i> , pages 45–54.	Luisa Ripoll Dominguez, Christian Widmer, Maxim-	1333
		ilian Wich, Hannah Danner, and Georg Groh. 2021.	1334
1279	Michael Gira, Ruisu Zhang, and Kangwook Lee. 2022.	End-to-end annotator bias approximation on crowd-	1335
1280	Debiasing pre-trained language models via efficient	sourced single-label sentiment analysis. In <i>Proceed-</i>	1336
1281	fine-tuning. In <i>Proceedings of the Second Workshop</i>	<i>ings of The Fourth International Conference on Natu-</i>	1337
1282	<i>on Language Technology for Equality, Diversity and</i>	<i>ral Language and Speech Processing (ICNLSP 2021)</i> ,	1338
1283	<i>Inclusion</i> , pages 59–69.	pages 1–10.	1339
1284	Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ri-	Felix Hamborg. 2020. Media bias, the social sciences,	1340
1285	cardo Muñoz Sánchez, Mugdha Pandya, and Adam	and nlp: automating frame analyses to identify bias	1341
1286	Lopez. 2021. Intrinsic bias metrics do not correlate	by word choice and labeling. In <i>Proceedings of the</i>	1342
1287	with application bias. In <i>Proceedings of the 59th An-</i>	<i>58th annual meeting of the association for computa-</i>	1343
1288	<i>annual Meeting of the Association for Computational</i>	<i>tional linguistics: student research workshop</i> , pages	1344
1289	<i>Linguistics and the 11th International Joint Confer-</i>	79–87.	1345
1290	<i>ence on Natural Language Processing (Volume 1:</i>		
1291	<i>Long Papers)</i> , pages 1926–1940.	Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021.	1346
		Decoupling adversarial training for fair nlp. In <i>Find-</i>	1347
1292	Hila Gonen and Yoav Goldberg. 2019. Lipstick on a	<i>ings of the Association for Computational Linguistics:</i>	1348
1293	pig: Debiasing methods cover up systematic gender	<i>ACL-IJCNLP 2021</i> , pages 471–477.	1349
1294	biases in word embeddings but do not remove them.		
1295	<i>arXiv preprint arXiv:1903.03862</i> .	Xudong Han, Timothy Baldwin, and Trevor Cohn. 2023.	1350
		Fair enough: Standardizing evaluation and model	1351
1296	Ana Valeria González, Maria Barrett, Rasmus Hvin-	selection for fairness research in nlp. In <i>Proceedings</i>	1352
1297	gelby, Kellie Webster, and Anders Søgaard. 2020.	<i>of the 17th Conference of the European Chapter of</i>	1353
1298	Type b reflexivization as an unambiguous testbed	<i>the Association for Computational Linguistics</i> , pages	1354
1299	for multilingual multi-task gender bias. In <i>Proceed-</i>	297–312.	1355
1300	<i>ings of the 2020 Conference on Empirical Methods</i>		
1301	<i>in Natural Language Processing (EMNLP)</i> , pages	Oussama Hansal, Ngoc Tan Le, and Fatiha Sadat. 2022.	1356
1302	2637–2648.	Indigenous language revitalization and the dilemma	1357
		of gender bias. In <i>Proceedings of the 4th Workshop</i>	1358
1303	Sian Gooding. 2022. On the ethical considerations of	<i>on Gender Bias in Natural Language Processing</i>	1359
1304	text simplification. In <i>Ninth Workshop on Speech</i>	<i>(GeBNLP)</i> , pages 244–254.	1360
1305	<i>and Language Processing for Assistive Technologies</i>		
1306	<i>(SLPAT-2022)</i> , pages 50–57.	Victor Petré Bach Hansen and Anders Søgaard. 2021.	1361
		Guideline bias in wizard-of-oz dialogues. In <i>Pro-</i>	1362
1307	Meiqi Guo, Rebecca Hwa, Yu-Ru Lin, and Wen-Ting	<i>ceedings of the 1st Workshop on Benchmarking: Past,</i>	1363
1308	Chung. 2020. Inflating topic relevance with ideology:	<i>Present and Future</i> , pages 8–14.	1364
1309	A case study of political ideology bias in social topic		

1365	Levon Haroutunian. 2022. Ethical considerations for low-resourced machine translation. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop</i> , pages 44–54.	1419
1366		1420
1367		1421
1368		1422
1369		1423
1370	Lucy Havens, Beatrice Alex, Benjamin Bach, and Melissa Terras. 2022. Uncertainty and inclusivity in gender bias annotation: An annotation taxonomy and annotated datasets of british english text. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 30–57.	1424
1371		
1372		
1373		
1374		
1375		
1376		
1377	Jacqueline He, Mengzhou Xia, Christiane Fellbaum, and Danqi Chen. 2022. Mabel: Attenuating gender bias using textual entailment data. <i>arXiv preprint arXiv:2210.14975</i> .	1425
1378		1426
1379		1427
1380		
1381	Zexue He, Bodhisattwa Prasad Majumder, and Julian McAuley. 2021. Detect and perturb: Neutral rewriting of biased and sensitive text via gradient-based decoding. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 4173–4181.	1428
1382		1429
1383		1430
1384		1431
1385		1432
1386		
1387	Brienna Herold, James Waller, and Raja Kushalnagar. 2022. Applying the stereotype content model to assess disability bias in popular pre-trained nlp models underlying ai-based assistive technologies. In <i>Ninth Workshop on Speech and Language Processing for Assistive Technologies (SLPAT-2022)</i> , pages 58–65.	1433
1388		1434
1389		1435
1390		1436
1391		1437
1392		1438
1393	Livnat Herzig, Alex Nunes, and Batia Snir. 2011. An annotation scheme for automated bias detection in wikipedia. In <i>Proceedings of the 5th Linguistic Annotation Workshop</i> , pages 47–55.	1439
1394		1440
1395		1441
1396		1442
1397	Marius Hessenthaler, Emma Strubell, Dirk Hovy, and Anne Lauscher. 2022. Bridging fairness and environmental sustainability in natural language processing. <i>arXiv preprint arXiv:2211.04256</i> .	1443
1398		1444
1399		
1400		
1401	Yasmeen Hitti, Eunbee Jang, Ines Moreno, and Carlyne Pelletier. 2019. Proposed taxonomy for gender bias in text; a filtering methodology for the gender generalization subtype. In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 8–17.	1445
1402		1446
1403		1447
1404		1448
1405		1449
1406		
1407	Dirk Hovy. 2015. Demographic factors improve classification performance. In <i>Proceedings of the 53rd annual meeting of the Association for Computational Linguistics and the 7th international joint conference on natural language processing (volume 1: Long papers)</i> , pages 752–762.	1450
1408		1451
1409		1452
1410		1453
1411		1454
1412		1455
1413	Dirk Hovy. 2018. The social and the neural network: How to make natural language processing about people again. In <i>Proceedings of the Second Workshop on Computational Modeling of People’s Opinions, Personality, and Emotions in Social Media</i> , pages 42–49.	1456
1414		1457
1415		1458
1416		1459
1417		1460
1418		1461
	Dirk Hovy, Federico Bianchi, and Tommaso Fornaciari. 2020. “you sound just like your father” commercial machine translation systems include stylistic biases. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1686–1690.	1462
		1463
		1464
		1465
	Dirk Hovy and Shrimai Prabhunoye. 2021. Five sources of bias in natural language processing. <i>Language and Linguistics Compass</i> , 15(8):e12432.	1466
		1467
		1468
		1469
	Dirk Hovy and Shannon L Spruit. 2016. The social impact of natural language processing. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 591–598.	1470
		1471
	Dirk Hovy and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice. In <i>The 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> . Association for Computational Linguistics.	1472
		1473
		1474
		1475
	Po-Sen Huang, Huan Zhang, Ray Jiang, Robert Stanforth, Johannes Welbl, Jack Rae, Vishal Maini, Dani Yogatama, and Pushmeet Kohli. 2020a. Reducing sentiment bias in language models via counterfactual evaluation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 65–83.	1476
		1477
		1478
		1479
		1480
	Tenghao Huang, Faeze Brahman, Vered Shwartz, and Snigdha Chaturvedi. 2021. Uncovering implicit gender bias in narratives through commonsense inference. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3866–3873.	1481
		1482
		1483
		1484
		1485
		1486
	Xiaolei Huang. 2022. Easy adaptation to mitigate gender bias in multilingual text classification. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 717–723.	1487
		1488
		1489
		1490
		1491
		1492
	Xiaolei Huang, Linzi Xing, Franck Dernoncourt, and Michael Paul. 2020b. Multilingual twitter corpus and baselines for evaluating demographic bias in hate speech recognition. In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 1440–1448.	1493
		1494
		1495
		1496
	Ben Hutchinson and Margaret Mitchell. 2019. 50 years of test (un) fairness: Lessons for machine learning. In <i>Proceedings of the conference on fairness, accountability, and transparency</i> , pages 49–58.	1497
		1498
		1499
		1500
	Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Social biases in nlp models as barriers for persons with disabilities. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5491–5501.	1501
		1502
		1503
		1504
		1505

1472	Sepehr Janghorbani and Gerard De Melo. 2023. Multi-modal bias: Introducing a framework for stereotypical bias assessment beyond gender and race in vision-language models. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1717–1727.	1528
1473		1529
1474		
1475		
1476		
1477		
1478	Sophie Jentzsch and Cigdem Turan. 2022. Gender bias in bert-measuring and analysing biases through sentiment rating in a realistic downstream classification task. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 184–199.	
1479		
1480		
1481		
1482		
1483		
1484	Shengyu Jia, Tao Meng, Jieyu Zhao, and Kai-Wei Chang. 2020. Mitigating gender bias amplification in distribution by posterior regularization. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2936–2942.	
1485		
1486		
1487		
1488		
1489	Meichun Jiao and Ziyang Luo. 2021. Gender bias hidden behind chinese word embeddings: The case of chinese adjectives. In <i>Proceedings of the 3rd Workshop on Gender Bias in Natural Language Processing</i> , pages 8–15.	
1490		
1491		
1492		
1493		
1494	Xisen Jin, Francesco Barbieri, Brendan Kennedy, Aida Mostafazadeh Davani, Leonardo Neves, and Xiang Ren. 2021a. On transferability of bias mitigation effects in language model fine-tuning. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3770–3783.	
1495		
1496		
1497		
1498		
1499		
1500		
1501		
1502	Zhijing Jin, Geeticka Chauhan, Brian Tse, Mrinmaya Sachan, and Rada Mihalcea. 2021b. How good is NLP? A sober look at NLP tasks through the lens of social impact. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3099–3113.	
1503		
1504		
1505		
1506		
1507		
1508	Przemyslaw Joniak and Akiko Aizawa. 2022. Gender biases and where to find them: Exploring gender bias in pre-trained transformer-based language models using movement pruning. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 67–73.	
1509		
1510		
1511		
1512		
1513		
1514	Lalitha Kameswari and Radhika Mamidi. 2021. Towards quantifying magnitude of political bias in news articles using a novel annotation schema. In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 671–678.	
1515		
1516		
1517		
1518		
1519		
1520	Lalitha Kameswari, Dama Sravani, and Radhika Mamidi. 2020. Enhancing bias detection in political news using pragmatic presupposition. In <i>Proceedings of the Eighth International Workshop on Natural Language Processing for Social Media</i> , pages 1–6.	
1521		
1522		
1523		
1524		
1525	Masahiro Kaneko and Danushka Bollegala. 2019. Gender-preserving debiasing for pre-trained word embeddings. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1641–1650.	1530
1526		1531
1527		1532
		1533
		1534
		1535
		1536
		1537
		1538
		1539
		1540
		1541
		1542
		1543
		1544
		1545
		1546
		1547
		1548
		1549
		1550
		1551
		1552
		1553
		1554
		1555
		1556
		1557
		1558
		1559
		1560
		1561
		1562
		1563
		1564
		1565
		1566
		1567
		1568
		1569
		1570
		1571
		1572
		1573
		1574
		1575
		1576
		1577
		1578
		1579
		1580
		1581
		1582
		1583

1584		<i>the Association for Computational Linguistics</i> , pages 2730–2743.	
1585			
1586	Vaibhav Kumar and Tenzin Bhotia. 2020. Fair embedding engine: A library for analyzing and mitigating gender bias in word embeddings. In <i>Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)</i> , pages 26–31.		
1587			
1588			
1589			
1590			
1591	Vaibhav Kumar, Tenzin Singhay Bhotia, and Tanmoy Chakraborty. 2020. Nurse is closer to woman than surgeon? mitigating gender-biased proximities in word embeddings. <i>Transactions of the Association for Computational Linguistics</i> , 8:486–503.		
1592			
1593			
1594			
1595			
1596	Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 166–172.		
1597			
1598			
1599			
1600			
1601	Alexander Kwako, Yixin Wan, Jieyu Zhao, Kai-Wei Chang, Li Cai, and Mark Hansen. 2022. Using item response theory to measure gender and racial bias of a bert-based automated english speech assessment system. In <i>Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)</i> , pages 1–7.		
1602			
1603			
1604			
1605			
1606			
1607			
1608	Faisal Ladhak, Esin Durmus, Mirac Suzgun, Tianyi Zhang, Dan Jurafsky, Kathleen Mckeown, and Tatsunori B Hashimoto. 2023. When do pre-training biases propagate to downstream tasks? a case study in text summarization. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3198–3211.		
1609			
1610			
1611			
1612			
1613			
1614			
1615	John P Lalor, Yi Yang, Kendall Smith, Nicole Forsgren, and Ahmed Abbasi. 2022. Benchmarking intersectional biases in NLP. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3598–3609.		
1616			
1617			
1618			
1619			
1620			
1621	Julian Lamont and Christi Favor. 2004. Distributive justice. <i>Handbook of political theory</i> , 1.		
1622			
1623	Brian N Larson. 2017. Gender as a variable in natural-language processing: Ethical considerations. In <i>First Workshop on Ethics in Natural Language Processing</i> . Association for Computational Linguistics.		
1624			
1625			
1626			
1627	Anne Lauscher and Goran Glavaš. 2019. Are we consistently biased? multidimensional analysis of biases in distributional word vectors. <i>NAACL HLT 2019</i> , page 85.		
1628			
1629			
1630			
1631	Anne Lauscher, Tobias Lueken, and Goran Glavaš. 2021. Sustainable modular debiasing of language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 4782–4797.		
1632			
1633			
1634			
1635	Konstantina Lazaridou, Alexander Löser, Maria Mestre, and Felix Naumann. 2020. Discovering biased news articles leveraging multiple human annotations. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 1268–1277.		
1636			
1637			
1638			
1639			
	Jochen L Leidner and Vassilis Plachouras. 2017. Ethical by design: Ethics best practices for natural language processing. In <i>Proceedings of the First ACL Workshop on Ethics in Natural Language Processing</i> , pages 30–40.		1640 1641 1642 1643 1644
	Heather Lent and Anders Søgaard. 2021. Common sense bias in semantic role labeling. In <i>Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT 2021)</i> , pages 114–119.		1645 1646 1647 1648
	Michael Lepori. 2020. Unequal representations: Analyzing intersectional biases in word embeddings using representational similarity analysis. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 1720–1728.		1649 1650 1651 1652 1653
	Shahar Levy, Koren Lazar, and Gabriel Stanovsky. 2021. Collecting a large-scale gender bias dataset for coreference resolution and machine translation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2470–2480.		1654 1655 1656 1657 1658
	Dave Lewis, Joss Moorkens, and Kaniz Fatema. 2017. Integrating the management of personal data protection and open science with research ethics. In <i>Proceedings of the First ACL Workshop on Ethics in Natural Language Processing</i> , pages 60–65.		1659 1660 1661 1662 1663
	Jiali Li, Shucheng Zhu, Ying Liu, and Pengyuan Liu. 2022a. Analysis of gender bias in social perception and judgement using chinese word embeddings. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 8–16.		1664 1665 1666 1667 1668 1669
	Tao Li, Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Vivek Srikumar. 2020. Uncovering stereotyping biases via underspecified questions. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3475–3489.		1670 1671 1672 1673 1674
	Yuantong Li, Xiaokai Wei, Zijian Wang, Shen Wang, Parminder Bhatia, Xiaofei Ma, and Andrew Arnold. 2022b. Debiasing neural retrieval via in-batch balancing regularization. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 58–66.		1675 1676 1677 1678 1679 1680
	Sheng Liang, Philipp Dufter, and Hinrich Schütze. 2020. Monolingual and multilingual reduction of gender bias in contextualized representations. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 5082–5093.		1681 1682 1683 1684 1685
	Sora Lim, Adam Jatowt, Michael Färber, and Masatoshi Yoshikawa. 2020. Annotating and analyzing biased sentences in news articles using crowdsourcing. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 1478–1484.		1686 1687 1688 1689 1690
	Tomasz Limisiewicz and David Mareček. 2022. Don’t forget about pronouns: Removing gender bias in language models without losing factual gender information. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 17–29.		1691 1692 1693 1694 1695 1696

1697	Emmy Liu, Michael Henry Tessler, Nicole Dubosh,	Nitin Madnani, Anastassia Loukina, Alina Von Davier,	1752
1698	Katherine Hiller, and Roger Levy. 2022a. Assessing	Jill Burstein, and Aoife Cahill. 2017. Building better	1753
1699	group-level gender bias in professional evaluations:	open-source tools to support fairness in automated	1754
1700	The case of medical student end-of-shift feedback.	scoring. In <i>Proceedings of the First ACL Workshop</i>	1755
1701	In <i>Proceedings of the 4th Workshop on Gender Bias</i>	<i>on Ethics in Natural Language Processing</i> , pages	1756
1702	in <i>Natural Language Processing (GeBNLP)</i> , pages	41–52.	1757
1703	86–93.		
1704	Haochen Liu, Jamell Dacon, Wenqi Fan, Hui Liu, Zitao	Vijit Malik, Sunipa Dev, Akihiro Nishi, Nanyun Peng,	1758
1705	Liu, and Jiliang Tang. 2020a. Does gender matter?	and Kai-Wei Chang. 2022. Socially aware bias mea-	1759
1706	towards fairness in dialogue systems. In <i>Proceedings</i>	surements for hindi language representations. In	1760
1707	<i>of the 28th International Conference on Computa-</i>	<i>Proceedings of the 2022 Conference of the North</i>	1761
1708	<i>tional Linguistics</i> , pages 4403–4416.	<i>American Chapter of the Association for Computa-</i>	1762
1709	Haochen Liu, Wei Jin, Hamid Karimi, Zitao Liu, and Jil-	<i>tional Linguistics: Human Language Technologies</i> ,	1763
1710	iang Tang. 2021. The authors matter: Understanding	pages 1041–1052.	1764
1711	and mitigating implicit bias in deep text classification.	Courtney Mansfield, Amandalynne Paullada, and Kris-	1765
1712	In <i>Findings of the Association for Computational Lin-</i>	ten Howell. 2022. Behind the mask: Demographic	1766
1713	<i>guistics: ACL-IJCNLP 2021</i> , pages 74–85.	bias in name detection for pii masking. In <i>Proceed-</i>	1767
1714	Haochen Liu, Joseph Thekinen, Sinem Mollaoglu,	<i>ings of the Second Workshop on Language Technol-</i>	1768
1715	Da Tang, Ji Yang, Youlong Cheng, Hui Liu, and	<i>ogy for Equality, Diversity and Inclusion</i> , pages 76–	1769
1716	Jiliang Tang. 2022b. Toward annotator group bias in	89.	1770
1717	crowdsourcing. In <i>Proceedings of the 60th Annual</i>	Thomas Manzini, Yao Chong Lim, Yulia Tsvetkov, and	1771
1718	<i>Meeting of the Association for Computational Lin-</i>	Alan W Black. 2019. Black is to criminal as cau-	1772
1719	<i>guistics (Volume 1: Long Papers)</i> , pages 1797–1806.	casian is to police: Detecting and removing multi-	1773
1720	Haochen Liu, Wentao Wang, Yiqi Wang, Hui Liu, Zitao	class bias in word embeddings. In <i>Proceedings of</i>	1774
1721	Liu, and Jiliang Tang. 2020b. Mitigating gender	<i>NAACL-HLT</i> , pages 615–621.	1775
1722	bias for neural dialogue generation with adversarial	Andrew Mao, Naveen Raman, Matthew Shu, Eric Li,	1776
1723	learning. In <i>Proceedings of the 2020 Conference on</i>	Franklin Yang, and Jordan Boyd-Graber. 2021. Elic-	1777
1724	<i>Empirical Methods in Natural Language Processing</i>	iting bias in question answering models through am-	1778
1725	<i>(EMNLP)</i> , pages 893–903.	biguity. In <i>Proceedings of the 3rd Workshop on Ma-</i>	1779
1726	Yunfei Long, Mingyu Ma, Qin Lu, Rong Xiang, and	<i>chine Reading for Question Answering</i> , pages 92–99.	1780
1727	Chu-Ren Huang. 2018. Dual memory network model	Sanjana Marcé and Adam Poliak. 2022. On gender	1781
1728	for biased product review classification. In <i>Pro-</i>	biases in offensive language classification models.	1782
1729	<i>ceedings of the 9th Workshop on Computational Ap-</i>	In <i>Proceedings of the 4th Workshop on Gender Bias</i>	1783
1730	<i>proaches to Subjectivity, Sentiment and Social Media</i>	<i>in Natural Language Processing (GeBNLP)</i> , pages	1784
1731	<i>Analysis</i> , pages 140–148.	174–183.	1785
1732	Anastassia Loukina, Nitin Madnani, and Klaus Zechner.	Sandra Martinková, Karolina Stanczak, and Isabelle	1786
1733	2019. The many dimensions of algorithmic fairness	Augenstein. 2023. Measuring gender bias in west	1787
1734	in educational applications. In <i>Proceedings of the</i>	slavic language models. In <i>Proceedings of the 9th</i>	1788
1735	<i>Fourteenth Workshop on Innovative Use of NLP for</i>	<i>Workshop on Slavic Natural Language Processing</i>	1789
1736	<i>Building Educational Applications</i> , pages 1–10.	<i>2023 (SlavicNLP 2023)</i> , pages 146–154.	1790
1737	Li Lucy and David Bamman. 2021. Gender and repre-	Judith Masthoff. 2011. Book review: Close engage-	1791
1738	sentation bias in gpt-3 generated stories. In <i>Proce-</i>	ments with artificial companions: Key social, psy-	1792
1739	<i>eedings of the Third Workshop on Narrative Under-</i>	chological, ethical, and design issues edited by yorick	1793
1740	<i>standing</i> , pages 48–55.	wilks. <i>Computational Linguistics</i> , 37(2).	1794
1741	Hongyin Luo and James Glass. 2023. Logic against	Camila M Mateo and David R Williams. 2020. Ad-	1795
1742	bias: Textual entailment mitigates stereotypical sen-	dressing bias and reducing discrimination: The pro-	1796
1743	tence reasoning. In <i>Proceedings of the 17th Confer-</i>	fessional responsibility of health care providers. <i>Aca-</i>	1797
1744	<i>ence of the European Chapter of the Association for</i>	<i>ademic Medicine</i> , 95(12S):S5–S10.	1798
1745	<i>Computational Linguistics</i> , pages 1235–1246.	Rowan Hall Maudslay, Hila Gonen, Ryan Cotterell, and	1799
1746	Xinyao Ma, Maarten Sap, Hannah Rashkin, and Yejin	Simone Teufel. 2019. It’s all in the name: Mitigating	1800
1747	Choi. 2020. Powertransformer: Unsupervised con-	gender bias with name-based counterfactual data sub-	1801
1748	trollable revision for biased language correction. In	stitution. In <i>Proceedings of the 2019 Conference on</i>	1802
1749	<i>Proceedings of the 2020 Conference on Empirical</i>	<i>Empirical Methods in Natural Language Processing</i>	1803
1750	<i>Methods in Natural Language Processing (EMNLP)</i> ,	<i>and the 9th International Joint Conference on Natu-</i>	1804
1751	pages 7426–7441.	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	1805
		5267–5275.	1806

1807	Chandler May, Alex Wang, Shikha Bordia, Samuel Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 622–628.	1865
1808		1866
1809		1867
1810		1868
1811		1869
1812		1870
1813		
1814	Elijah Mayfield, Michael Madaio, Shrimai Prabhumoye, David Gerritsen, Brittany McLaughlin, Ezekiel Dixon-Román, and Alan W Black. 2019. Equity beyond bias in language technologies for education. In <i>Proceedings of the Fourteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 444–460.	1871
1815		1872
1816		1873
1817		
1818		1874
1819		1875
1820		1876
		1877
1821	Nicholas Meade, Elinor Poole-Dayán, and Siva Reddy. 2022. An empirical survey of the effectiveness of debiasing techniques for pre-trained language models. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1878–1898.	1878
1822		1879
1823		1880
1824		1881
1825		1882
1826		
1827	Michal Měchura. 2022. A taxonomy of bias-causing ambiguities in machine translation. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 168–173.	1883
1828		1884
1829		1885
1830		1886
		1887
1831	Ninareh Mehrabi, Umang Gupta, Fred Morstatter, Greg Ver Steeg, and Aram Galstyan. 2022. Attributing fair decisions with attention interventions. In <i>Proceedings of the 2nd Workshop on Trustworthy Natural Language Processing (TrustNLP 2022)</i> , pages 12–25.	1888
1832		1889
1833		1890
1834		1891
1835		1892
1836		
1837	Jack Merullo, Luke Yeh, Abram Handler, Alvin Grisom II, Brendan O’Connor, and Mohit Iyyer. 2019. Investigating sports commentator bias within a large corpus of american football broadcasts. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6355–6361.	1893
1838		1894
1839		1895
1840		1896
1841		1897
1842		1898
1843		1899
1844		
1845	Josh Meyer, Lindy Rauchenstein, Joshua D Eisenberg, and Nicholas Howell. 2020. Artie bias corpus: An open dataset for detecting demographic bias in speech applications. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 6462–6468.	1900
1846		1901
1847		1902
1848		1903
1849		1904
1850		1905
1851	Piotr Miłkowski, Marcin Gruza, Kamil Kanclerz, Przemysław Kazienko, Damian Grimling, and Jan Kočoń. 2021. Personal bias in prediction of emotions elicited by textual opinions. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing</i> , pages 248–259.	1906
1852		1907
1853		1908
1854		1909
1855		1910
1856		1911
1857		1912
1858		
1859	Fatemehsadat Miresheghallah and Taylor Berg-Kirkpatrick. 2021. Style pooling: Automatic text style obfuscation for improved classification fairness. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 2009–2022.	1913
1860		1914
1861		1915
1862		1916
1863		
1864		1917
		1918
		1919
		1920
		1921
	Pushkar Mishra, Helen Yannakoudakis, and Ekaterina Shutova. 2021. Modeling users and online communities for abuse detection: A position on ethics and explainability. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3374–3385.	
	Saif Mohammad. 2022a. Ethics sheet for automatic emotion recognition and sentiment analysis. <i>Computational Linguistics</i> , 48(2):239–278.	
	Saif Mohammad. 2022b. Ethics sheets for AI tasks. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8368–8379.	
	Rodrigo Alejandro Chávez Mulsa and Gerasimos Spanakis. 2020. Evaluating bias in dutch word embeddings. In <i>Proceedings of the Second Workshop on Gender Bias in Natural Language Processing</i> , pages 56–71.	
	Robert Munro and Alex Carmen Morrison. 2020. Detecting independent pronoun bias with partially-synthetic data generation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2011–2017.	
	Taichi Murayama, Shoko Wakamiya, and Eiji Aramaki. 2021. Mitigation of diachronic bias in fake news detection dataset. In <i>Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT 2021)</i> , pages 182–188.	
	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. Stereoset: Measuring stereotypical bias in pretrained language models. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5356–5371.	
	Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1953–1967.	
	Aurélie Névéol, Yoann Dupont, Julien Bezançon, and Karèn Fort. 2022. French crows-pairs: Extending a challenge dataset for measuring social bias in masked language models to a language other than english. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8521–8531.	
	Malvina Nissim, Rik van Noord, and Rob van der Goot. 2020. Fair is better than sensational: Man is to doctor as woman is to doctor. <i>Computational Linguistics</i> , 46(2):487–497.	
	Timothy Niven and Hung-Yu Kao. 2020. Measuring alignment to authoritarian state media as framing bias. In <i>Proceedings of the 3rd NLP4IF Workshop on NLP for Internet Freedom: Censorship, Disinformation, and Propaganda</i> , pages 11–21.	

1922	Debora Nozza, Federico Bianchi, and Dirk Hovy. 2022.	Shrimai Prabhumoye, Brendon Boldt, Ruslan Salakhutdinov, and Alan W Black. 2021.	1977
1923	Pipelines for social bias testing of large language	Case study: Deontological ethics in NLP. In <i>Proceedings of the 2021</i>	1978
1924	models. In <i>Proceedings of BigScience Episode# 5–</i>	<i>Conference of the North American Chapter of the</i>	1979
1925	<i>Workshop on Challenges & Perspectives in Creating</i>	<i>Association for Computational Linguistics: Human</i>	1980
1926	<i>Large Language Models</i> , pages 68–74.	<i>Language Technologies</i> , pages 3784–3798.	1981
1927	Hadas Orgad and Yonatan Belinkov. 2022. Choose your	Flavien Prost, Nithum Thain, and Tolga Bolukbasi. 2019.	1983
1928	lenses: Flaws in gender bias evaluation . In <i>Proceed-</i>	Debiasing embeddings for reduced gender	1984
1929	<i>ings of the 4th Workshop on Gender Bias in Natu-</i>	bias in text classification. In <i>Proceedings of the First</i>	1985
1930	<i>ral Language Processing (GeBNLP)</i> , pages 151–167,	<i>Workshop on Gender Bias in Natural Language Pro-</i>	1986
1931	Seattle, Washington. Association for Computational	<i>cessing</i> , pages 69–75.	1987
1932	Linguistics.	Ming Qian, Jessie Liu, Chaofeng Li, and Liming Pals. 2019a.	1988
1933	Hadas Orgad, Seraphina Goldfarb-Tarrant, and Yonatan	2019a. A comparative study of english-chinese trans-	1989
1934	Belinkov. 2022. How gender debiasing affects in-	lations of court texts by machine and human trans-	1990
1935	ternal model representations, and why it matters. In	lators and the word2vec based similarity measure’s	1991
1936	<i>Proceedings of the 2022 Conference of the North</i>	ability to gauge human evaluation biases. In <i>Proceed-</i>	1992
1937	<i>American Chapter of the Association for Computa-</i>	<i>ings of Machine Translation Summit XVII: Translator,</i>	1993
1938	<i>tional Linguistics: Human Language Technologies</i> ,	<i>Project and User Tracks</i> , pages 95–100.	1994
1939	pages 2602–2628.	Rebecca Qian, Candace Ross, Jude Fernandes, Eric	1995
1940	Martin Orr, Kirsten Van Kessel, and Dave Parry. 2022.	Smith, Douwe Kiela, and Adina Williams. 2022.	1996
1941	The ethical role of computational linguistics in digi-	Perturbation augmentation for fairer nlp. <i>arXiv preprint</i>	1997
1942	tal psychological formulation and suicide prevention.	<i>arXiv:2205.12586</i> .	1998
1943	In <i>Proceedings of the Eighth Workshop on Computa-</i>	Yusu Qian, Urwa Muaz, Ben Zhang, and Jae Won Hyun. 2019b.	1999
1944	<i>tional Linguistics and Clinical Psychology</i> , pages	2019b. Reducing gender bias in word-level language	2000
1945	17–29.	models with a gender-equalizing loss function. In	2001
1946	Prasanna Parasurama and João Sedoc. 2022. Gendered	<i>Proceedings of the 57th Annual Meeting of the Asso-</i>	2002
1947	language in resumes and its implications for algo-	<i>ciation for Computational Linguistics: Student Re-</i>	2003
1948	rithmic bias in hiring. In <i>Proceedings of the 4th</i>	<i>search Workshop</i> , pages 223–228.	2004
1949	<i>Workshop on Gender Bias in Natural Language Pro-</i>	Krithika Ramesh, Gauri Gupta, and Sanjay Singh. 2021.	2005
1950	<i>cessing (GeBNLP)</i> , pages 74–74.	Evaluating gender bias in hindi-english machine	2006
1951	Ji Ho Park, Jamin Shin, and Pascale Fung. 2018. Re-	translation. In <i>Proceedings of the 3rd Workshop on</i>	2007
1952	ducing gender bias in abusive language detection.	<i>Gender Bias in Natural Language Processing</i> , pages	2008
1953	In <i>Proceedings of the 2018 Conference on Empiri-</i>	16–23.	2009
1954	<i>cal Methods in Natural Language Processing</i> , pages	Krithika Ramesh, Sunayana Sitaram, and Monojit	2010
1955	2799–2804.	Choudhury. 2023. Fairness in language models be-	2011
1956	Mihir Parmar, Swaroop Mishra, Mor Geva, and Chitta	yond english: Gaps and challenges. In <i>Findings</i>	2012
1957	Baral. 2022. Don’t blame the annotator: Bias already	<i>of the Association for Computational Linguistics:</i>	2013
1958	starts in the annotation instructions. <i>arXiv preprint</i>	<i>EACL 2023</i> , pages 2061–2074.	2014
1959	<i>arXiv:2205.00415</i> .	Alan Ramponi and Sara Tonelli. 2022. Features or spu-	2015
1960	Alicia Parrish, Angelica Chen, Nikita Nangia,	rious artifacts? data-centric baselines for fair and	2016
1961	Vishakh Padmakumar, Jason Phang, Jana Thompson,	robust hate speech detection. In <i>Proceedings of the</i>	2017
1962	Phu Mon Htut, and Samuel Bowman. 2022. Bbq: A	<i>2022 Conference of the North American Chapter of</i>	2018
1963	hand-built bias benchmark for question answering.	<i>the Association for Computational Linguistics: Hu-</i>	2019
1964	In <i>Findings of the Association for Computational</i>	<i>man Language Technologies</i> , pages 3027–3040. As-	2020
1965	<i>Linguistics: ACL 2022</i> , pages 2086–2105.	sociation for Computational Linguistics.	2021
1966	Jannik Pedersen, Martin Laursen, Pernille Vinholt,	Adithya Renduchintala, Denise Díaz, Kenneth Heafield,	2022
1967	Anne Alnor, and Thiusius Savarimuthu. 2023. In-	Xian Li, and Mona Diab. 2021. Gender bias ampli-	2023
1968	vestigating anatomical bias in clinical machine learn-	fication during speed-quality optimization in neural	2024
1969	ing algorithms. In <i>Findings of the Association for</i>	machine translation. In <i>Proceedings of the 59th An-</i>	2025
1970	<i>Computational Linguistics: EACL 2023</i> , pages 1368–	<i>nual Meeting of the Association for Computational</i>	2026
1971	1380.	<i>Linguistics and the 11th International Joint Confer-</i>	2027
1972	Matúš Pikuliak, Ivana Beňová, and Viktor Bachratý.	<i>ence on Natural Language Processing (Volume 2:</i>	2028
1973	2023. In-depth look at word filling societal bias mea-	<i>Short Papers)</i> , pages 99–109.	2029
1974	asures. In <i>Proceedings of the 17th Conference of the</i>	Alisa Rieger, Mariët Theune, and Nava Tintarev. 2020.	2030
1975	<i>European Chapter of the Association for Computa-</i>	Toward natural language mitigation strategies for cog-	2031
1976	<i>tional Linguistics</i> , pages 3630–3647.	nitive biases in recommender systems. In <i>2nd Work-</i>	2032
		<i>shop on Interactive Natural Language Technology</i>	2033
		<i>for Explainable Artificial Intelligence</i> , pages 50–54.	2034

2035	Anthony Rios, Reenam Joshi, and Hejin Shin. 2020.	<i>Linguistics: Human Language Technologies</i> , pages	2091
2036	Quantifying 60 years of gender bias in biomedical re-	5884–5906, Seattle, United States. Association for	2092
2037	search with word embeddings. In <i>Proceedings of the</i>	Computational Linguistics.	2093
2038	<i>19th SIGBioMed Workshop on Biomedical Language</i>		
2039	<i>Processing</i> .		
2040	Alexey Romanov, Maria De-Arteaga, Hanna Wal-	Danielle Saunders and Bill Byrne. 2020. Reducing gen-	2094
2041	lach, Jennifer Chayes, Christian Borgs, Alexandra	bias in neural machine translation as a domain	2095
2042	Chouldechova, Sahin Geyik, Krishnaram Kenthapadi,	adaptation problem. In <i>Proceedings of the 58th An-</i>	2096
2043	Anna Rumshisky, and Adam Kalai. 2019. What’s	<i>annual Meeting of the Association for Computational</i>	2097
2044	in a name? reducing bias in bios without access	<i>Linguistics</i> , pages 7724–7736.	2098
2045	to protected attributes. In <i>Proceedings of the 2019</i>		
2046	<i>Conference of the North American Chapter of the</i>	Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Mat-	2099
2047	<i>Association for Computational Linguistics: Human</i>	teo Negri, and Marco Turchi. 2021. Gender bias in	2100
2048	<i>Language Technologies, Volume 1 (Long and Short</i>	machine translation. <i>Transactions of the Association</i>	2101
2049	<i>Papers)</i> , pages 4187–4195.	<i>for Computational Linguistics</i> , 9:845–874.	2102
2050	Candace Ross, Boris Katz, and Andrei Barbu. 2021.	Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Mat-	2103
2051	Measuring social biases in grounded vision and lan-	teo Negri, and Marco Turchi. 2022. Under the mor-	2104
2052	guage embeddings. In <i>Proceedings of the 2021 Con-</i>	phosyntactic lens: A multifaceted evaluation of gen-	2105
2053	<i>ference of the North American Chapter of the Asso-</i>	der bias in speech translation. In <i>Proceedings of the</i>	2106
2054	<i>ciation for Computational Linguistics: Human Lan-</i>	<i>60th Annual Meeting of the Association for Compu-</i>	2107
2055	<i>guage Technologies</i> , pages 998–1008.	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	2108
2056	Rachel Rudinger, Chandler May, and Benjamin	1807–1824.	2109
2057	Van Durme. 2017. Social bias in elicited natural		
2058	language inferences. In <i>Proceedings of the First ACL</i>	Ramit Sawhney, Arshiya Aggarwal, and Rajiv Shah.	2110
2059	<i>Workshop on Ethics in Natural Language Processing</i> ,	2021. An empirical investigation of bias in the mul-	2111
2060	pages 74–79.	timodal analysis of financial earnings calls. In <i>Pro-</i>	2112
2061	Rachel Rudinger, Jason Naradowsky, Brian Leonard,	<i>ceedings of the 2021 Conference of the North Amer-</i>	2113
2062	and Benjamin Van Durme. 2018. Gender bias in	<i>ican Chapter of the Association for Computational</i>	2114
2063	coreference resolution. In <i>Proceedings of NAACL-</i>	<i>Linguistics: Human Language Technologies</i> , pages	2115
2064	<i>HLT</i> , pages 8–14.	3751–3757.	2116
2065	Magnus Sahlgren and Fredrik Olsson. 2019. Gender	Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021.	2117
2066	bias in pretrained swedish embeddings. In <i>Proceed-</i>	Self-diagnosis and self-debiasing: A proposal for re-	2118
2067	<i>ings of the 22nd Nordic Conference on computational</i>	ducing corpus-based bias in nlp. <i>Transactions of the</i>	2119
2068	<i>linguistics</i> , pages 35–43.	<i>Association for Computational Linguistics</i> , 9:1408–	2120
2069	Sebastin Santy, Anku Rani, and Monojit Choudhury.	1424.	2121
2070	2021. Use of formal ethical reviews in NLP literature:	Natalie Schluter. 2018. The glass ceiling in NLP. In	2122
2071	Historical trends and current practices. In <i>Findings of</i>	<i>Proceedings of the 2018 Conference on Empirical</i>	2123
2072	<i>the Association for Computational Linguistics: ACL-</i>	<i>Methods in Natural Language Processing</i> , pages	2124
2073	<i>IJCNLP 2021</i> , pages 4704–4710.	2793–2798.	2125
2074	Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi,	Katja Geertruida Schmahl, Tom Julian Viering, Stavros	2126
2075	and Noah A Smith. 2019. The risk of racial bias in	Makrodimitris, Arman Naseri Jahfari, David Tax, and	2127
2076	hate speech detection. In <i>Proceedings of the 57th</i>	Marco Loog. 2020. Is wikipedia succeeding in re-	2128
2077	<i>annual meeting of the association for computational</i>	ducing gender bias? assessing changes in gender bias	2129
2078	<i>linguistics</i> , pages 1668–1678.	in wikipedia using word embeddings. In <i>Proceed-</i>	2130
2079	Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Juraf-	<i>ings of the Fourth Workshop on Natural Language</i>	2131
2080	sky, Noah A Smith, and Yejin Choi. 2020. Social	<i>Processing and Computational Social Science</i> , pages	2132
2081	bias frames: Reasoning about social and power im-	94–103.	2133
2082	plications of language. In <i>Proceedings of the 58th</i>	Tyler Schnoebelen. 2017. Goal-oriented design for ethi-	2134
2083	<i>Annual Meeting of the Association for Computational</i>	cal machine learning and NLP. In <i>Proceedings of the</i>	2135
2084	<i>Linguistics</i> , pages 5477–5490.	<i>First ACL Workshop on Ethics in Natural Language</i>	2136
2085	Maarten Sap, Swabha Swayamdipta, Laura Vianna,	<i>Processing</i> , pages 88–93.	2137
2086	Xuhui Zhou, Yejin Choi, and Noah A. Smith. 2022.	Stephanie Schoch, Diyi Yang, and Yangfeng Ji. 2020.	2138
2087	Annotators with attitudes: How annotator beliefs	“This is a problem, don’t you agree?” Framing and	2139
2088	and identities bias toxic language detection . In <i>Pro-</i>	bias in human evaluation for natural language gener-	2140
2089	<i>ceedings of the 2022 Conference of the North Amer-</i>	ation. In <i>Proceedings of the 1st Workshop on Evalu-</i>	2141
2090	<i>ican Chapter of the Association for Computational</i>	<i>ating NLG Evaluation</i> , pages 10–16.	2142
		Lena Schwertmann, Manoj Prabhakar Kannan Ravi,	2143
		and Gerard De Melo. 2023. Model-agnostic bias	2144
		measurement in link prediction. In <i>Findings of the</i>	2145
		<i>Association for Computational Linguistics: EACL</i>	2146
		<i>2023</i> , pages 1587–1603.	2147

2148	João Sedoc and Lyle Ungar. 2019. The role of protected class word lists in bias identification of contextualized word representations. In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 55–61.	2205
2149		2206
2150		2207
2151		2208
2152		2209
2153	Indira Sen, Mattia Samory, Fabian Flöck, Claudia Wagner, and Isabelle Augenstein. 2021. How does counterfactually augmented data impact models for social computing constructs? In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 325–344.	2210
2154		2211
2155		2212
2156		2213
2157		2214
2158		2215
2159	Indira Sen, Mattia Samory, Claudia Wagner, and Isabelle Augenstein. 2022. Counterfactually augmented data and unintended bias: The case of sexism and hate speech detection. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4716–4726.	2216
2160		2217
2161		2218
2162		2219
2163		2220
2164		2221
2165		2222
2166	Emeralda Sesari, Max Hort, and Federica Sarro. 2022. An empirical study on the fairness of pre-trained word embeddings. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 129–144.	2223
2167		2224
2168		2225
2169		2226
2170		2227
2171	Usman Shahid, Barbara Di Eugenio, Andrew Rojecki, and Elena Zheleva. 2020. Detecting and understanding moral biases in news. In <i>Proceedings of the First Joint Workshop on Narrative Understanding, Storylines, and Events</i> , pages 120–125.	2228
2172		2229
2173		2230
2174		2231
2175		2232
2176	Aili Shen, Xudong Han, Trevor Cohn, Timothy Baldwin, and Lea Frermann. 2022. Optimising equal opportunity fairness in model training. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4073–4084.	2233
2177		2234
2178		2235
2179		2236
2180		2237
2181		2238
2182	Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3407–3412.	2239
2183		2240
2184		2241
2185		2242
2186		2243
2187		2244
2188		2245
2189	Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2020. Towards controllable biases in language generation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3239–3254.	2246
2190		2247
2191		2248
2192		2249
2193		2250
2194	Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2021. Societal biases in language generation: Progress and challenges. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4275–4293.	2251
2195		2252
2196		2253
2197		2254
2198		2255
2199		2256
2200		2257
2201	Emily Sheng and David C Uthus. 2020. Investigating societal biases in a poetry composition system. In <i>Proceedings of the Second Workshop on Gender Bias in Natural Language Processing</i> , pages 93–106.	2258
2202		2259
2203		2260
2204		2260
	Seungjae Shin, Kyungwoo Song, JoonHo Jang, Hyemi Kim, Weonyoung Joo, and Il-Chul Moon. 2020. Neutralizing gender bias in word embeddings with latent disentanglement and counterfactual generation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3126–3140.	2205
		2206
		2207
		2208
		2209
		2210
	Boaz Shmueli, Jan Fell, Soumya Ray, and Lun-Wei Ku. 2021. Beyond fair pay: Ethical implications of NLP crowdsourcing. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3758–3769.	2211
		2212
		2213
		2214
		2215
		2216
	Vered Shwartz and Yejin Choi. 2020. Do neural language models overcome reporting bias? In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6863–6870.	2217
		2218
		2219
		2220
	Andrew Silva, Pradyumna Tambwekar, and Matthew Gombolay. 2021. Towards a comprehensive understanding and accurate evaluation of societal biases in pre-trained transformers. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2383–2389.	2221
		2222
		2223
		2224
		2225
		2226
		2227
	Charese Smiley, Frank Schilder, Vassilis Plachouras, and Jochen L Leidner. 2017. Say the right thing right: Ethics issues in natural language generation systems. In <i>Proceedings of the First ACL Workshop on Ethics in Natural Language Processing</i> , pages 103–108.	2228
		2229
		2230
		2231
		2232
		2233
	Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. 2022. “i’m sorry to hear that”: Finding new biases in language models with a holistic descriptor dataset. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9180–9211.	2234
		2235
		2236
		2237
		2238
		2239
	Timo Spinde, Manuel Plank, Jan-David Krieger, Terry Ruas, Bela Gipp, and Akiko Aizawa. 2021. Neural media bias detection using distant supervision with babe-bias annotations by experts. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1166–1177.	2240
		2241
		2242
		2243
		2244
		2245
	Maximilian Spliethöver and Henning Wachsmuth. 2020. Argument from old man’s view: Assessing social bias in argumentation. In <i>Proceedings of the 7th Workshop on Argument Mining</i> , pages 76–87.	2246
		2247
		2248
		2249
	Tejas Srinivasan and Yonatan Bisk. 2022a. Worst of both worlds: Biases compound in pre-trained vision-and-language models. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 77–85, Seattle, Washington. Association for Computational Linguistics.	2250
		2251
		2252
		2253
		2254
		2255
	Tejas Srinivasan and Yonatan Bisk. 2022b. Worst of both worlds: Biases compound in pre-trained vision-and-language models. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 77–85.	2256
		2257
		2258
		2259
		2260

2261	Artūrs Stefanovičs, Toms Bergmanis, and Mārcis Pinnis.	Chris Sweeney and Maryam Najafian.	2315
2262	2020. Mitigating gender bias in machine translation	2019. A transparent framework for evaluating unintended demo-	2316
2263	with target gender annotations. In <i>Proceedings of</i>	graphic bias in word embeddings. In <i>Proceedings</i>	2317
2264	<i>of the Fifth Conference on Machine Translation</i> , pages	<i>of the 57th Annual Meeting of the Association for</i>	2318
2265	629–638.	<i>Computational Linguistics</i> , pages 1662–1667.	2319
2266	Gabriel Stanovsky, Noah A Smith, and Luke Zettle-	Masashi Takeshita, Yuki Katsumata, Rafal Rzepka, and	2320
2267	moyer. 2019. Evaluating gender bias in machine	Kenji Araki. 2020. Can existing methods debias lan-	2321
2268	translation. In <i>Proceedings of the 57th Annual Meet-</i>	guages other than english? first attempt to analyze	2322
2269	<i>ing of the Association for Computational Linguistics</i> ,	and mitigate japanese word embeddings. In <i>Pro-</i>	2323
2270	pages 1679–1684.	<i>ceedings of the Second Workshop on Gender Bias in</i>	2324
2271	Ryan Steed, Swetasudha Panda, Ari Kobren, and	<i>Natural Language Processing</i> , pages 44–55.	2325
2272	Michael Wick. 2022. Upstream mitigation is not	Yarden Tal, Inbal Magar, and Roy Schwartz. 2022.	2326
2273	all you need: Testing the bias transfer hypothesis in	Fewer errors, but more stereotypes? the effect of	2327
2274	pre-trained language models. In <i>Proceedings of the</i>	model size on gender bias. In <i>Proceedings of the</i>	2328
2275	<i>60th Annual Meeting of the Association for Compu-</i>	<i>4th Workshop on Gender Bias in Natural Language</i>	2329
2276	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	<i>Processing (GeBNLP)</i> , pages 112–120.	2330
2277	3524–3542.	Zeerak Talat, Hagen Blix, Josef Valvoda, Maya Indira	2331
2278	Angus Stevenson. 2010. <i>Oxford dictionary of English</i> .	Ganesh, Ryan Cotterell, and Adina Williams. 2022a.	2332
2279	Oxford University Press, USA.	On the machine learning of ethical judgments from	2333
2280	Shivashankar Subramanian, Xudong Han, Timothy	natural language. In <i>Proceedings of the 2022 Con-</i>	2334
2281	Baldwin, Trevor Cohn, and Lea Frermann. 2021a.	<i>ference of the North American Chapter of the Asso-</i>	2335
2282	Evaluating debiasing techniques for intersectional	<i>ciation for Computational Linguistics: Human Lan-</i>	2336
2283	biases. In <i>Proceedings of the 2021 Conference on</i>	<i>guage Technologies</i> , pages 769–779.	2337
2284	<i>Empirical Methods in Natural Language Processing</i> ,	Zeerak Talat, Aurelie Neveol, Stella Biderman, Miruna	2338
2285	pages 2492–2498.	Cliniciu, Manan Dey, Shayne Longpre, Sasha Luc-	2339
2286	Shivashankar Subramanian, Afshin Rahimi, Timothy	cioni, Maraim Masoud, Margaret Mitchell, Dragomir	2340
2287	Baldwin, Trevor Cohn, and Lea Frermann. 2021b.	Radev, et al. 2022b. You reap what you sow: On	2341
2288	Fairness-aware class imbalanced learning. In <i>Pro-</i>	the challenges of bias evaluation under multilingual	2342
2289	<i>ceedings of the 2021 Conference on Empirical Meth-</i>	settings. In <i>Proceedings of BigScience Episode# 5–</i>	2343
2290	<i>ods in Natural Language Processing</i> , pages 2045–	<i>Workshop on Challenges & Perspectives in Creating</i>	2344
2291	2051.	<i>Large Language Models</i> , pages 26–41.	2345
2292	Jiao Sun and Nanyun Peng. 2021. Men are elected,	Rachael Tatman. 2017a. Gender and dialect bias in	2346
2293	women are married: Events gender bias on wikipedia.	youtube’s automatic captions. In <i>Proceedings of the</i>	2347
2294	In <i>Proceedings of the 59th Annual Meeting of the</i>	<i>first ACL workshop on ethics in natural language</i>	2348
2295	<i>Association for Computational Linguistics and the</i>	<i>processing</i> , pages 53–59.	2349
2296	<i>11th International Joint Conference on Natural Lan-</i>	Rachael Tatman. 2017b. “oh, i’ve heard that before”:	2350
2297	<i>guage Processing (Volume 2: Short Papers)</i> , pages	Modelling own-dialect bias after perceptual learning	2351
2298	350–360.	by weighting training data. In <i>Proceedings of the 7th</i>	2352
2299	Tianxiang Sun, Junliang He, Xipeng Qiu, and Xuanjing	<i>Workshop on Cognitive Modeling and Computational</i>	2353
2300	Huang. 2022. Bertscore is unfair: On social bias	<i>Linguistics (CMCL 2017)</i> , pages 29–34.	2354
2301	in language model-based metrics for text generation.	Wenyi Tay. 2019. Not all reviews are equal: Towards ad-	2355
2302	<i>arXiv preprint arXiv:2210.07626</i> .	ressing reviewer biases for opinion summarization.	2356
2303	Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang,	In <i>Proceedings of the 57th Annual Meeting of the</i>	2357
2304	Mai ElSherief, Jieyu Zhao, Diba Mirza, Eliza-	<i>Association for Computational Linguistics: Student</i>	2358
2305	beth Belding, Kai-Wei Chang, and William Yang	<i>Research Workshop</i> , pages 34–42.	2359
2306	Wang. 2019. Mitigating gender bias in natural lan-	Andamlak Terkik, Emily Prud’hommeaux, Cecilia	2360
2307	guage processing: Literature review. <i>arXiv preprint</i>	Ovesdotter Alm, Christopher Homan, and Scott	2361
2308	<i>arXiv:1906.08976</i> .	Franklin. 2016. Analyzing gender bias in student	2362
2309	Simon Šuster, Stéphan Tulkens, and Walter Daelemans.	evaluations. In <i>Proceedings of COLING 2016, the</i>	2363
2310	2017. A short review of ethical challenges in clinical	<i>26th International Conference on Computational Lin-</i>	2364
2311	natural language processing. In <i>Proceedings of the</i>	<i>guistics: Technical Papers</i> , pages 868–876.	2365
2312	<i>First ACL Workshop on Ethics in Natural Language</i>	Ewoenam Tokpo, Pieter Delobelle, Bettina Berendt, and	2366
2313	<i>Processing</i> , pages 80–87, Valencia, Spain. Associa-	Toon Calders. 2023. How far can it go? on intrinsic	2367
2314	tion for Computational Linguistics.	gender bias mitigation for text classification. In	2368
		<i>Proceedings of the 17th Conference of the European</i>	2369
		<i>Chapter of the Association for Computational Lin-</i>	2370
		<i>guistics</i> , pages 3410–3425.	2371

2372	Ewoenam Kwaku Tokpo and Toon Calders. 2022. Text style transfer for bias mitigation using masked language modeling. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop</i> , pages 163–171.	2427
2373		2428
2374		2429
2375		2430
2376		2431
2377		2432
2378		2433
2379	Autumn Toney and Aylin Caliskan. 2021. Valnorm quantifies semantics to reveal consistent valence biases across languages and over centuries. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7203–7218.	2434
2380		2435
2381		2436
2382		2437
2383		2438
2384		2439
2385	Samia Touileb, Lilja Øvrelid, and Erik Velldal. 2021. Using gender-and polarity-informed models to investigate bias. In <i>Proceedings of the 3rd Workshop on Gender Bias in Natural Language Processing</i> , pages 66–74.	2440
2386		2441
2387		2442
2388		2443
2389	Samia Touileb, Lilja Øvrelid, and Erik Velldal. 2022. Occupational biases in norwegian and multilingual language models. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 200–211.	2444
2390		2445
2391		2446
2392		2447
2393		2448
2394	Samia Touileb, Lilja Øvrelid, and Erik Velldal. 2023. Measuring normative and descriptive biases in language models using census data. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 2234–2240.	2449
2395		2450
2396		2451
2397		2452
2398		2453
2399		2454
2400	Jonas Troles and Ute Schmid. 2021. Extending challenge sets to uncover gender bias in machine translation impact of stereotypical verbs and adjectives. <i>WMT 2021</i> , page 531.	2455
2401		2456
2402		2457
2403		2458
2404	Dimitrios Tsarapatsanis and Nikolaos Aletras. 2021. On the ethical limits of natural language processing on legal text. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3590–3599.	2459
2405		2460
2406		2461
2407		2462
2408		2463
2409	Oskar Van Der Wal, Jaap Jumelet, Katrin Schulz, and Willem Zuidema. 2022. The birth of bias: A case study on the evolution of gender bias in an english language model. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 75–75.	2464
2410		2465
2411		2466
2412		2467
2413		2468
2414		2469
2415	Lindsey Vanderlyn, Gianna Weber, Michael Neumann, Dirk Vāth, Sarina Meyer, and Ngoc Thang Vu. 2021. “it seemed like an annoying woman”: On the perception and ethical considerations of affective language in text-based conversational agents. In <i>Proceedings of the 25th Conference on Computational Natural Language Learning</i> , pages 44–57.	2470
2416		2471
2417		2472
2418		2473
2419		2474
2420		2475
2421		2476
2422	Francisco Vargas and Ryan Cotterell. 2020. Exploring the linear subspace hypothesis in gender bias mitigation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2902–2913.	2477
2423		2478
2424		2479
2425		2480
2426		2481
	Eva Maria Vecchi, Neele Falk, Iman Jundi, and Gabriella Lapesa. 2021. Towards argument mining for social good: A survey. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1338–1352.	2482
	Pranav Narayanan Venkit, Sanjana Gautam, Ruchi Panchanadikar, Ting-Hao Huang, and Shomir Wilson. 2023. Nationality bias in text generation. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 116–122.	
	Xuan-Son Vu, Thanh-Son Nguyen, Duc-Trong Le, and Lili Jiang. 2020. Multimodal review generation with privacy and fairness awareness. In <i>28th International Conference on Computational Linguistics (COLING), Barcelona, Spain (Online), December 8-13, 2020.</i> , pages 414–425. International Committee on Computational LinguisticsInternational Committee	
	Bo Wang, Tao Shen, Guodong Long, Tianyi Zhou, and Yi Chang. 2021a. Eliminating sentiment bias for aspect-level sentiment classification with unsupervised opinion extraction. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3002–3012.	
	Jialu Wang, Yang Liu, and Xin Wang. 2021b. Are gender-neutral queries really gender-neutral? mitigating gender bias in image search. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 1995–2008.	
	Jianan Wang, Xin Wang, Fang Li, Zhen Xu, Zhuoran Wang, and Baoxun Wang. 2017. Group linguistic bias aware neural response generation. In <i>Proceedings of the 9th SIGHAN Workshop on Chinese Language Processing</i> , pages 1–10.	
	Jun Wang, Benjamin Rubinstein, and Trevor Cohn. 2022. Measuring and mitigating name biases in neural machine translation. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2576–2590.	
	Liwen Wang, Yuanmeng Yan, Keqing He, Yanan Wu, and Weiran Xu. 2021c. Dynamically disentangling social bias from task-oriented representations with adversarial attack. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3740–3750.	
	Shuo Wang, Zhaopeng Tu, Zhixing Tan, Shuming Shi, Maosong Sun, and Yang Liu. 2021d. On the language coverage bias for neural machine translation. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4778–4790.	
	Tianlu Wang, Xi Victoria Lin, Nazneen Fatema Rajani, Bryan McCann, Vicente Ordonez, and Caiming	

2483	Xiong. 2020. Double-hard debias: Tailoring word embeddings for gender bias mitigation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5443–5453.	
2484		
2485		
2486		
2487	Zhao Wang, Kai Shu, and Aron Culotta. 2021e. Enhancing model robustness and fairness with causality: A regularization approach. In <i>Proceedings of the First Workshop on Causal Inference and NLP</i> , pages 33–43.	
2488		
2489		
2490		
2491		
2492	Zeerak Waseem, Smarika Lulz, Joachim Bingel, and Isabelle Augenstein. 2021. Disembodied machine learning: On the illusion of objectivity in NLP. <i>arXiv preprint arXiv:2101.11974</i> .	
2493		
2494		
2495		
2496	Melvin Wevers. 2019. Using word embeddings to examine gender bias in dutch newspapers, 1950-1990. In <i>Proceedings of the 1st International Workshop on Computational Approaches to Historical Language Change</i> , pages 92–97.	
2497		
2498		
2499		
2500		
2501	Maximilian Wich, Hala Al Kuwatly, and Georg Groh. 2020a. Investigating annotator bias with a graph-based approach. In <i>Proceedings of the Fourth Workshop on Online Abuse and Harms</i> , pages 191–199.	
2502		
2503		
2504		
2505	Maximilian Wich, Jan Bauer, and Georg Groh. 2020b. Impact of politically biased data on hate speech classification. In <i>Proceedings of the fourth workshop on online abuse and harms</i> , pages 54–64.	
2506		
2507		
2508		
2509	Maximilian Wich, Christian Widmer, Gerhard Hagerer, and Georg Groh. 2021. Investigating annotator bias in abusive language datasets. In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 1515–1525.	
2510		
2511		
2512		
2513		
2514		
2515	Williams Institute. 2021. 1.2 million LGBTQ adults in the US identify as nonbinary .	
2516		
2517	Robert Wolfe and Aylin Caliskan. 2021. Low frequency names exhibit bias and overfitting in contextualizing language models. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 518–532.	
2518		
2519		
2520		
2521		
2522	Mengzhou Xia, Anjalie Field, and Yulia Tsvetkov. 2020. Demoting racial bias in hate speech detection. In <i>Proceedings of the Eighth International Workshop on Natural Language Processing for Social Media</i> , pages 7–14.	
2523		
2524		
2525		
2526		
2527	Zhongbin Xie, Vid Kocijan, Thomas Lukasiewicz, and Oana-Maria Camburu. 2023. Counter-gap: Counterfactual bias evaluation through gendered ambiguous pronouns. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3743–3755.	
2528		
2529		
2530		
2531		
2532		
2533	Tae Yano, Philip Resnik, and Noah A Smith. 2010. Shedding (a thousand points of) light on biased language. In <i>Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk</i> , pages 152–158.	
2534		
2535		
2536		
2537		
	Samira Zad, Joshuan Jimenez, and Mark Finlayson. 2021. Hell hath no fury? correcting bias in the nrc emotion lexicon. In <i>Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021)</i> , pages 102–113.	2538 2539 2540 2541 2542
	Haiyang Zhang, Alison Sneyd, and Mark Stevenson. 2020. Robustness and reliability of gender bias assessment in word embeddings: The role of base pairs. In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing</i> , pages 759–769.	2543 2544 2545 2546 2547 2548 2549 2550
	Jieyu Zhao and Kai-Wei Chang. 2020. Logan: Local group bias detection by clustering. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1968–1977.	2551 2552 2553 2554 2555
	Jieyu Zhao, Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Kai-Wei Chang. 2021. Ethical-advice taker: Do language models understand natural language interventions? In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4158–4164.	2556 2557 2558 2559 2560 2561
	Jieyu Zhao, Subhabrata Mukherjee, Kai-Wei Chang, Ahmed Hassan Awadallah, et al. 2020. Gender bias in multilingual embeddings and cross-lingual transfer. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> .	2562 2563 2564 2565 2566
	Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2979–2989.	2567 2568 2569 2570 2571 2572
	Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 15–20.	2573 2574 2575 2576 2577 2578 2579
	Ruiqi Zhong, Yanda Chen, Desmond Patton, Charlotte Selous, and Kathleen Mckeown. 2019. Detecting and reducing bias in a high stakes domain. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4765–4775.	2580 2581 2582 2583 2584 2585 2586
	Yang Zhong, Jingfeng Yang, Wei Xu, and Diyi Yang. 2021. Wikibias: Detecting multi-span subjective biases in language. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1799–1814.	2587 2588 2589 2590 2591
	Pei Zhou, Weijia Shi, Jieyu Zhao, Kuan-Hao Huang, Muhao Chen, Ryan Cotterell, and Kai-Wei Chang.	2592 2593

2594 2019. Examining gender bias in languages with
2595 grammatical gender. In *Proceedings of the 2019 Con-*
2596 *ference on Empirical Methods in Natural Language*
2597 *Processing and the 9th International Joint Confer-*
2598 *ence on Natural Language Processing (EMNLP-*
2599 *IJCNLP)*, pages 5276–5284.

2600 Xuhui Zhou, Maarten Sap, Swabha Swayamdipta, Yejin
2601 Choi, and Noah A Smith. 2021. Challenges in au-
2602 tomated debiasing for toxic language detection. In
2603 *Proceedings of the 16th Conference of the European*
2604 *Chapter of the Association for Computational Lin-*
2605 *guistics: Main Volume*, pages 3143–3155.

2606 Yi Zhou, Masahiro Kaneko, and Danushka Bollegala.
2607 2022. Sense embeddings are also biased—evaluating
2608 social biases in static and contextualised sense em-
2609 beddings. In *Proceedings of the 60th Annual Meet-*
2610 *ing of the Association for Computational Linguistics*
2611 *(Volume 1: Long Papers)*, pages 1924–1935.

2612 Caleb Ziems, Jane Yu, Yi-Chia Wang, Alon Halevy,
2613 and Diyi Yang. 2022. The moral integrity corpus: A
2614 benchmark for ethical dialogue systems. In *Proceed-*
2615 *ings of the 60th Annual Meeting of the Association for*
2616 *Computational Linguistics (Volume 1: Long Papers)*,
2617 pages 3755–3773.

2618 Nadezhda Zueva, Madina Kabirova, and Pavel Kalaidin.
2619 2020. Reducing unintended identity bias in russian
2620 hate speech detection. In *Proceedings of the Fourth*
2621 *Workshop on Online Abuse and Harms*, pages 65–69.

7 Appendix

Type	Papers
Survey/Position Paper	Drugan and Babych (2010), Masthoff (2011), Escartín et al. (2016), Hovy and Spruit (2016), Leidner and Plachouras (2017), et al. (2017), Schnoebelen (2017), Larson (2017), De Jong et al. (2018), Mayfield et al. (2019), Sun et al. (2019), Blodgett et al. (2021), Bender et al. (2020), Schoch et al. (2020), Savoldi et al. (2021), Mishra et al. (2021), Czarnowska et al. (2021), Dougherty et al. (2021), Sheng et al. (2021), Hovy and Yang (2021), Santy et al. (2021), (2022b), Chandrabose et al. (2021), Vecchi et al. (2021), Talat et al. (2022), Balkir et al. (2022), Orgad and Belinkov (2022), Akyürek et al. (2022), Mohammad (2022a), Orr et al. (2022), Benotti and Blodgett et al. (2023), Benotti et al. (2023), Ramesh et al. (2023), Rios et al. (2023)
Dataset, Benchmark	Yano et al. (2010), Flekova et al. (2016), Braun et al. (2016), Kallmeyer et al. (2016), Rudinger et al. (2018), Zhao et al. (2018), Kiritchenko and Mooney (2019), Gaut et al. (2020), Friedman et al. (2019), Hitti et al. (2019), Fan et al. (2019), Liu et al. (2020a), Sap et al. (2020), Lou et al. (2021), Lim et al. (2020), Shahid et al. (2020), Lazaridou et al. (2020), Uthus (2020), Gala et al. (2020), Dinan et al. (2020b), Meyer et al. (2020b), Gangula et al. (2019), Zueva et al. (2020), Nangia et al. (2020), Blodgett et al. (2021), Spinde et al. (2021), Zhong et al. (2021), Søgaard (2021), Zhao et al. (2021), Aksenov et al. (2021), Kamnitsas et al. (2021), Renduchintala et al. (2021), Barikeri et al. (2021), Miłkowski et al. (2021), et al. (2022), Ziems et al. (2022), Zhou et al. (2022), Akyürek et al. (2022), Bansal et al. (2022), Smith et al. (2022), Qian et al. (2022), et al. (2023), Madnani et al. (2017), Kumar and Bhotia (2020), (2021), Zad et al. (2021), Nozza et al. (2022), Měchura (2022)
Methods	Gonen and Goldberg (2019), Wevers (2019), Hovy et al. (2020), Nissim et al. (2020), Wang et al. (2021d), Wolfe and Caliskan (2021), et al. (2021), Herold et al. (2022), Bertsch et al. (2022), Rampal et al. (2022), Cao et al. (2022a), Alshahrani et al. (2022), Daems and Hackenbusch et al. (2015), Terkik et al. (2016), Tatman (2017b), Wang et al. (2017a), Tatman (2017a), Zhao et al. (2017), Long et al. (2018), Baly et al. (2019), Kloppenburg (2019), Färber et al. (2019), Zhong et al. (2019), Tetreault et al. (2019), and Ungar (2019), Gangula et al. (2019), Karve et al. (2019), Iqbal and Najafian (2019), Chaloner and Maldonado (2019), Qian et al. (2019), Kaneko and Bollegala (2019), Huang et al. (2020a), Bhat et al. (2020), Dinan et al. (2020a), Bordia and Bowman (2019), Manzini et al. (2020), et al. (2020), Schwertmann et al. (2023), Pedersen et al. (2023), Bauer et al. (2023), Kumar et al. (2023), Touileb et al. (2023), Parmar et al. (2022), et al. (2023), Alemany et al. (2023), Martinková et al. (2023), Iqbal et al. (2022), Cao et al. (2022b), Sap et al. (2020), Xia et al. (2020), et al. (2022), Orgad et al. (2022), Sen et al. (2022), Gupta et al. (2022), Gira et al. (2022), Guo et al. (2022), Wang et al. (2022), An et al. (2022), Borchers et al. (2022), Touileb et al. (2022), Jentzsch and et al. (2022), Tal et al. (2022), Li et al. (2022b), Limisiewicz et al. (2022), Ahn et al. (2022), Joniak and Aizawa (2022), Chen et al. (2022), Du et al. (2022), Agrawal et al. (2022), Sesari et al. (2022), Steed et al. (2022), Schick et al. (2021), Zhou et al. (2021), Kane et al. (2021), Lucy and Bamman (2021), Jin et al. (2021a), Iqbal et al. (2021b), Silva et al. (2021), Sun and Peng (2021), Wang et al. (2021), Han et al. (2021), Sen et al. (2021), Ramesh et al. (2021), Subramanian and Oh (2021), Davidson et al. (2019), May et al. (2019), Sap et al. (2020), Piper (2020), Guo et al. (2020), Vu et al. (2020), Spliethöver et al. (2020), Munro and Morrison (2020), Zhao and Chang (2020), Schmahla et al. (2020), Chen et al. (2020b), Mulsa and Spanakis (2020), Shiran et al. (2020), Cotterell (2020), Liang et al. (2020), Wich et al. (2020a), Jia et al. (2020), Dev et al. (2021a), Kumar et al. (2020), Field and Tsvetkov et al. (2020a), Stafanovičs et al. (2020), Li et al. (2020), Ma et al. (2020), Al Kuwatly et al. (2020), Bartl et al. (2020), Liu et al. (2020b), Farkas et al. (2020), Sheng et al. (2020), Gaci et al. (2022), Hutchinson et al. (2020), et al. (2021e), Van Der Wal et al. (2022), Hansen and Søgaard (2022), Azarpanah and Farhadloo (2021), Jiao and Luo (2021), Touileb et al. (2021), Ciora et al. (2021), Murayama et al. (2021), Gaido et al. (2021), et al. (2022), Dawkins (2021a), Gillis (2021), Miresghallah and Bhat et al. (2021), Wang et al. (2021a), Malik et al. (2022), Subramanian et al. (2022), Garimella et al. (2021), Ghosh et al. (2021), Mehrabi et al. (2022)
Both above	Goldfarb-Tarrant et al. (2021), Janghorbani and De Melo (2022), Troles and Schmid (2021), Liu et al. (2022a), Névéol et al. (2022)
Analysis	Herzig et al. (2011), Rudinger et al. (2017), Lauscher and Glavač et al. (2021), Vanderlyn et al. (2021), Falenska and Çetinoğlu (2021), Sap et al. (2021), Bhat et al. (2021), Haroutunian (2022), Qian et al. (2021)